

## Deep Neural Networks for Inverse Problems with Pseudodifferential Operators: An Application to Limited-Angle Tomography\*

Tatiana A. Bubba<sup>†</sup>, Mathilde Galinier<sup>‡</sup>, Matti Lassas<sup>†</sup>, Marco Prato<sup>‡</sup>, Luca Ratti<sup>†</sup>, and  
Samuli Siltanen<sup>†</sup>

**Abstract.** We propose a novel convolutional neural network (CNN), called  $\Psi$ DONet, designed for learning pseudodifferential operators ( $\Psi$ DOs) in the context of linear inverse problems. Our starting point is the iterative soft thresholding algorithm (ISTA), a well-known algorithm to solve sparsity-promoting minimization problems. We show that, under rather general assumptions on the forward operator, the unfolded iterations of ISTA can be interpreted as the successive layers of a CNN, which in turn provides fairly general network architectures that, for a specific choice of the parameters involved, allow us to reproduce ISTA, or a perturbation of ISTA for which we can bound the coefficients of the filters. Our case study is the limited-angle X-ray transform and its application to limited-angle computed tomography (LA-CT). In particular, we prove that, in the case of LA-CT, the operations of upscaling, downscaling, and convolution, which characterize our  $\Psi$ DONet and most deep learning schemes, can be exactly determined by combining the convolutional nature of the limited-angle X-ray transform and basic properties defining an orthogonal wavelet system. We test two different implementations of  $\Psi$ DONet on simulated data from limited-angle geometry, generated from the ellipse data set. Both implementations provide equally good and noteworthy preliminary results, showing the potential of the approach we propose and paving the way to applying the same idea to other convolutional operators which are  $\Psi$ DOs or Fourier integral operators.

**Key words.** x-ray transform, limited-angle tomography, deep neural networks, convolutional neural networks, wavelets, sparse regularization, Fourier integral operators, pseudodifferential operators, microlocal analysis

**AMS subject classifications.** 44A12, 68T07, 35S30, 58J40, 92C55

**DOI.** 10.1137/20M1343075

**1. Introduction.** In the context of microlocal analysis, the theory of pseudodifferential operators ( $\Psi$ DOs), introduced by Kohn and Nirenberg in 1965, and Fourier integral operators (FIOs), defined by Hörmander in 1971, finds remarkable applications in many fields of

\*Received by the editors June 5, 2020; accepted for publication (in revised form) October 26, 2020; published electronically May 3, 2021.

<https://doi.org/10.1137/20M1343075>

**Funding:** The work of the first, third, fifth, and sixth authors was supported by the Finnish Centre of Excellence in Inverse Modelling and Imaging, 2018–2025, grant 312339, the Academy of Finland grants 284715, 312110, 310822, and the Faculty of Science ATMATH project. The work of the second and fourth authors was supported by the INdAM-GNCS research projects 2018. The work of the second author was supported by the INdAM Doctoral Programme in Mathematics and/or Applications cofunded by Marie Skłodowska-Curie Actions (INdAM-DP-COFUND-2015) grant 713485. The work of the fourth author was partially supported by the ECSEL JU programme under the PRYSTINE Project grant 783190.

<sup>†</sup>Department of Mathematics and Statistics, University of Helsinki, Helsinki, FI-00014 ([tatiana.bubba@helsinki.fi](mailto:tatiana.bubba@helsinki.fi), [matti.lassas@helsinki.fi](mailto:matti.lassas@helsinki.fi), [luca.ratti@helsinki.fi](mailto:luca.ratti@helsinki.fi), [samuli.siltanen@helsinki.fi](mailto:samuli.siltanen@helsinki.fi)).

<sup>‡</sup>Department of Physics, Computer Science and Mathematics, University of Modena and Reggio Emilia, Modena, IT-41125 ([mathildeemmanuelle.galinier@unimore.it](mailto:mathildeemmanuelle.galinier@unimore.it), [marco.prato@unimore.it](mailto:marco.prato@unimore.it)).

mathematics, from spectral theory to general relativity, from the study of the behavior of chaotic systems to scattering theory, and inverse problems [26, 27]. A prominent example in the inverse problem field is given by the X-ray transform or, in the two-dimensional case, Radon transform:

$$(1.1) \quad R(u)(s, \omega) = \int_{-\infty}^{\infty} u(s\omega^\perp + t\omega) dt, \quad s \in \mathbb{R}, \omega, \omega^\perp \in S^1,$$

where  $\omega^\perp$  denotes the vector in the unit sphere  $S^1$  obtained by rotating  $\omega$  counterclockwise by  $90^\circ$  [39]. It is possible to show (see, e.g., [44]) that the normal operator  $R^*R$  of the Radon transform  $R$  is an elliptic  $\Psi$ DO of order  $-1$  and a convolutional operator associated with the Calderón–Zygmund kernel  $K(x, y) = \frac{1}{|x-y|}$  for  $x \neq y$ . When the direction vector  $\omega$  is restricted within a limited angular range  $[-\Gamma, \Gamma]$ , the normal operator  $R_\Gamma^*R_\Gamma$  of the limited-angle Radon transform  $R_\Gamma$  is a convolutional operator associated with the kernel

$$(1.2) \quad K(x, y) = \frac{1}{|x-y|} \chi_\Gamma(x-y) \quad \text{for } x \neq y,$$

where  $\chi_\Gamma$  denotes the indicator function of the cone in  $\mathbb{R}^2$  between the angles  $-\Gamma$  and  $\Gamma$ . The operator  $R_\Gamma^*R_\Gamma$  is no longer a  $\Psi$ DO, but it belongs to a wider class of FIOs, which includes operators associated with a kernel showing some discontinuities along lines [27].

The inverse problem arising from the limited angle Radon transform, i.e., limited-angle computed tomography (LA-CT), appears frequently in practical applications, such as dental tomography [32], damage detection in concrete structures [24], breast tomosynthesis [56], or electron tomography [15].

Microlocal analysis has been widely applied on the Radon transform, especially in the case of incomplete data, with the purpose to characterize its behavior with respect to singularities in the images (see, e.g., [20, 45, 46, 30, 29, 16]). In particular, some recent works are focused on the treatment of artifacts appearing in reconstructions from limited-angle data (see [40, 41, 7]). In this framework, microlocal analysis can be used to predict which singularities can be reconstructed in a stable way from limited-angle measurements [7, 18, 33, 43]. In practice, thanks to microlocal analysis, we are able to identify the part of the wavefront set of the target corresponding to the missing wedge from the measurement geometry.

Even with this fundamental information, the task of robustly recovering the unknown quantity of interest from such partial indirect measurement is a challenging one, due to the ill-posedness of the CT problem, which is even more severe because of the limited angular range [13]. As a result, classical methods, such as the filtered backprojection (FBP) [39], yield poor performances. Traditional inversion methods of the form (2.3)–(2.5), based on complementing the insufficient measurements by imposing a priori information on the solution, define effective regularization methods which generally allow for accurate reconstructions from fewer tomographic measurements than usually required by standard methods like FBP. In more recent years, machine learning approaches, in particular, deep learning, with convolutional neural networks (CNNs) being the most prominent design in the context of imaging, are increasingly impacting the field of inverse problems [4], and (LA-)CT is no exception (see, in particular, [4, section 4] for an overview of learning approaches from a functional analytic

regularization perspective and [4, section 7] for their applicability to prototypical examples of inverse problems, including CT). The majority of recent data-driven approaches for LA-CT focuses on recovering or inpainting the missing part of the wavefront set from the measured data (see, e.g., [9, 49] and the references therein for a thorough review of model-based and data-driven approaches based on sparsifying transforms and edge-preserving regularizers in the context of LA-CT).

In this paper, we are *not* interested in designing an(other) approach for inferring the missing wedge in LA-CT, but rather we aim at investigating neural networks inspired by FIOs and  $\Psi$ DOs, for which LA-CT is a case study. Our starting point is the traditional sparsity-based minimization problem of the form (2.3)–(2.5). A well-known technique for its solution is the iterative soft thresholding algorithm (ISTA), introduced in 2004 in the seminal paper by Daubechies, Defrise, and De Mol [12]. The convergence result in the paper relies on the assumption that the sparsifying system forms an orthogonal basis, as is the case for many families of wavelets [36]. ISTA iteratively creates the sequence  $\{w^{(n)}\}_{n=1}^N$  as follows:

$$(1.3) \quad w^{(n)} = \mathcal{S}_{\lambda/L} \left( w^{(n-1)} - \frac{1}{L} K^{(n)} w^{(n-1)} + \frac{1}{L} b^{(n)} \right),$$

where, in our case,  $K^{(n)} = W R_{\Gamma}^* R_{\Gamma} W^*$  with  $W$  wavelet transform associated with an orthogonal family,  $b^{(n)} = W R_{\Gamma}^* m$ , and  $\mathcal{S}_{\beta}(w)$  is the (componentwise) soft-thresholding operator (see (2.7) and (3.1) for all the details). It is well known that the unrolled iterations of ISTA can be considered as the layers of a neural network. Learned ISTA (LISTA), introduced in [22], and ISTA-Net, introduced in [54], are examples of neural networks obtained by laying out the operations of ISTA for a few iterations. The major difference with our approach is that LISTA and ISTA-Net are not CNNs. Unrolled schemes coming from proximal primal-dual optimization methods are also proposed in [2, 3], where the proximal operators are replaced with CNNs. While in [2, 3] the goal is to learn a proximal operator, in the approach we propose the regularization term is fixed and we learn a correction of the normal operator  $R_{\Gamma}^* R_{\Gamma}$ . Deep unfolded schemes for problems other than CT are introduced, for instance, in [25, 55]. In [35] the authors propose an unsupervised approach, combined with unrolled schemes, to learn adversarial regularizers and apply it to the case of full-angle CT. In [28] the authors investigate the relationship between CNNs and iterative optimization methods, including ISTA, for the case of normal operators associated with a forward model which is a convolution. However, the resulting U-net, FBPCnvNet, does not aim at imitating an unrolled version of an iterative method, which makes it fundamentally different in spirit to the methodology we propose. Indeed, the goal of our work is to show that, under some assumptions on the operator  $R_{\Gamma} W^*$ , it is possible to interpret the operations in (1.3) as a layer of a CNN, which in turn provides fairly general network architectures that allow us to recover standard ISTA for a specific choice of the parameters involved.

Motivated by this, we propose a new CNN, which we name  $\Psi$ DONet, aimed at learning convolutional FIOs and  $\Psi$ DOs. The key feature of  $\Psi$ DONet is that we split the convolutional kernel into  $K = K_0 + K_1$ , where  $K_0$  is the known part of the model (in the limited-angle case,  $K_0 = R_{\Gamma}^* R_{\Gamma}$ ) and  $K_1$  is an unknown  $\Psi$ DO to be determined or, better, to be learned. Basically, in  $K_1$  lays the potential to add information in the reconstruction process with respect to the known part of the model  $K_0$ .  $\Psi$ DONet takes advantage of the possibility to use small

filters encoding a combination of upscaling, downscaling, and convolution operations, as it is common practice in deep learning. Remarkably, we prove that such operations can be exactly determined combining the convolutional nature of the limited-angle Radon transform and basic properties defining an orthogonal wavelet system. While this might seem contrary to the machine learning philosophy which finds its strength in avoiding any predefined structure for neural networks, our recipe gives insight into understanding and interpreting the results of the proposed CNN, combining results from FIOs,  $\Psi$ DOs, and classical variational regularization theory. At the same time, the possibility to deploy such operations allows for a significant reduction of the parameters involved, especially when compared to the standard interpretation of ISTA as a recurrent neural network: this is fundamental when it comes to a practical numerical implementation of the proposed CNN. Overall,  $\Psi$ DONet is able to reproduce ISTA, or a perturbation of ISTA for which we can bound the coefficients of the filters, and has the potential to learn  $\Psi$ DO-like structures which are intrinsic to the problem at hand.

As a proof of concept, we test  $\Psi$ DONet on simulated data from limited-angle geometry, generated from the ellipse data set. We provide two different implementations of  $\Psi$ DONet: filter-based  $\Psi$ DONet ( $\Psi$ DONet-F), where the backprojection operator is approximated by its filter equivalent, and operator-based  $\Psi$ DONet ( $\Psi$ DONet-O), where the backprojection operator encoded in  $K_0$  is not approximated but explicitly computed. Both implementations provide equally good and noteworthy preliminary results, the main difference being a greater computational efficiency for  $\Psi$ DONet-O. The improvement provided by our results, compared to standard ISTA (and FBP), bodes well for further numerical testing which we leave to future work.

Finally, we stress that the contribution of our paper is mainly theoretical and is in line with current research in data-driven inversion, which combines knowledge from traditional inverse problems theory with data-driven techniques. While in our paper we derived the result contingently to the case of limited-angle Radon transform, our approach is actually very general and can be extended to any convolutional operator which is a FIO or  $\Psi$ DO. This is the case, for instance, of the geodesic X-ray transform [50], and its applications in seismic imaging, or synthetic-aperture radar [42]. Finally, our paper paves the way to theoretical generalization results, in light of recent contributions like [14].

The remainder of this paper is organized as follows: [section 2](#) is devoted to reviewing the theoretical background of sparsity promoting regularization, and the wavelet transform. In [section 3](#), we detail the key idea of our approach, namely, we give a convolutional interpretation of ISTA using the wavelet transform. The neural network architecture we propose,  $\Psi$ DONet, is introduced in [section 4](#), where we also prove our main theoretical result. Two different implementations of  $\Psi$ DONet, which we call filter-based  $\Psi$ DONet and operator-based  $\Psi$ DONet, are described in [section 5](#). Finally, we demonstrate the performance of our network by a series of numerical experiments (see [section 6](#)). Concluding remarks and future prospects are briefly summarized in [section 7](#). The appendices collect proofs of some of the results presented in [section 2](#).

**2. Theoretical background.** In this section, we collect some theoretical results which are preliminary to the main discussion of the paper.

**2.1. Sparsity-promoting regularization via ISTA.** Consider the inverse problem of determining  $u^\dagger \in X$  from the measurements  $m = Au^\dagger + \epsilon$ , being  $A : X \rightarrow Y$  a linear bounded operator between the Hilbert spaces  $X$  and  $Y$ . The perturbation  $\epsilon \in Y$  is such that  $\|\epsilon\|_Y \leq \delta$ . The main application we have in mind is the limited-angle Radon transform  $R_\Gamma$ , which is a continuous linear operator, e.g., from  $X = L^2(\Omega)$  (being  $\Omega \subset \mathbb{R}^2$ ) to  $Y = L^2([-\Gamma, \Gamma] \times [-S, S])$  (see [39, Theorem 2.10]).

Introduce an orthonormal basis  $\{\psi_I\}_{I \in \mathbb{N}}$  in  $X$ . For later purposes, we will assume that such a basis is a wavelet system. Define  $W : X \rightarrow \ell^2(\mathbb{N})$  as the operator associating with any  $u \in X$  the sequence of its components with respect to the wavelet basis  $(Wu)_I = (u, \psi_I)_X$ , where  $(\cdot, \cdot)_X$  denotes the inner product in  $X$ . We assume to know a priori that the exact solution  $u^\dagger$  is sparse with respect to the wavelet basis  $\psi_I$ :

$$(2.1) \quad Wu^\dagger = w^\dagger \in \ell^0(\mathbb{N}).$$

The reconstruction of  $u^\dagger$  (or, equivalently,  $w^\dagger$ ) from the noisy measurements  $m$  is in general an ill-posed problem, hence we introduce the following regularized problem:

$$(2.2) \quad \min_{w \in \ell^1(\mathbb{N})} \|AW^*w - m\|_Y^2 + \lambda \|w\|_{\ell^1}$$

with  $\lambda > 0$ . The requirement  $w \in \ell^1(\mathbb{N})$  is in general not satisfied by any  $w = Wu$ ,  $u \in X$ ; hence, we define  $Z \subset X$ ,  $Z = \{u \in X : Wu \in \ell^1(\mathbb{N})\}$ . In particular, in the tomography application, it is possible to show that the  $\ell^1$  norm of the components of the wavelet representation of an  $L^2(\Omega)$  function is equivalent to the Besov norm  $B_{1,1}^1(\Omega)$  (see, e.g., [12, formula (A3)]). Hence, the minimization problem (2.2) is equivalent to

$$(2.3) \quad \min_{u \in Z} \|Au - m\|_Y^2 + \lambda \|u\|_Z.$$

It is well known that the regularization term involving the  $\ell^1$  norm is a good choice to encode the a priori information regarding the sparsity of  $w^\dagger$ . In particular if the noise level tends to 0, there exists a suitable choice of  $\lambda = \lambda(\delta)$  ensuring the convergence of  $w_\lambda^\delta$  to  $w^\dagger$  with  $w_\lambda^\delta$  the solution of (2.2). We report a result from [17] which also shows that such a convergence occurs with linear rate. In particular, [17, Corollary 2] does not require  $w^\dagger$  to satisfy a classical source condition, but relies on the sparsity assumption (2.1) and on the injectivity of the operator  $A$ . Such a property can be restrictive in some applications, and as a consequence many alternative results involve some weaker assumptions (as the well-known restricted isometry property); nevertheless, in our tomographic application, we can rely on the injectivity of the Radon transform, even in the limited-angle case.

**Proposition 2.1.** *Let  $w^\dagger$  satisfy (2.1), and suppose  $A : X \rightarrow Y$  is injective. Define  $w_\delta^\lambda$  as a solution of problem (2.2) associated with a regularization parameter  $\lambda$  and a noise level  $\delta$ . For sufficiently small  $\delta$ , provided that  $\lambda$  is chosen such that  $\lambda = c_0\delta$ , then there exists a positive constant  $c_1 = c_1(c_0, A, \|w^\dagger\|_{\ell^0})$  such that*

$$(2.4) \quad \|w^\dagger - w_\delta^\lambda\|_{\ell^1} \leq c_1\delta.$$

This proposition is an immediate consequence of [17, Corollary 2], relying on [17, Lemma 2] to ensure that  $A$  is weak\*-to-weak continuous. From now on, we suppose that  $\lambda$  is chosen as a linear function of  $\delta$  and denote  $u_\delta^\lambda$  as  $u_\delta$  and  $w_\delta^\lambda$  as  $w_\delta$ .

We now introduce a finite-dimensional approximation of the regularized problem (2.2). Consider the subspace  $X_p \subset X$ ,  $X_p = \text{span}\{\psi_I\}_{I=1}^p$ , mapped by  $W$  into the space  $W_p = \{w \in \mathbb{R}^N : w_I = 0 \ \forall I > p\}$  (which is isomorphic to  $\mathbb{R}^p$ ). Denote by  $\mathbb{P}_p$  the orthogonal projection of  $\ell^2(\mathbb{N})$  onto  $W_p$  and by  $\tilde{\mathbb{P}}_p = W^* \mathbb{P}_p W$  the orthogonal projection of  $X$  onto  $X_p$ . Moreover, we introduce an orthogonal basis  $\{\varphi_j\}_{j=1}^\infty$  on  $Y$  and define  $Y_q = \text{span}\{\varphi_j\}_{j=1}^q$  and the projection  $\mathbb{P}_q : Y \rightarrow Y_q$ . For any choice of  $p, q > 0$ , let  $A_{p,q}$  be the representation of the operator  $A$  in the subspaces  $X_p, Y_q$ , namely,  $A_{p,q} = \mathbb{P}_q A \tilde{\mathbb{P}}_p^*$ . Consider the following minimization problem:

$$(2.5) \quad \min_{w \in W_p} \|A_{p,q} W^* w - \mathbb{P}_q m\|_Y^2 + \lambda \|w\|_{\ell^1}.$$

Denote by  $w_{\delta,p,q}$  a solution of (2.5). We can prove the following convergence result.

**Proposition 2.2.** *Let  $w^\dagger$  satisfy (2.1) and  $A$  be an injective operator. Suppose moreover that for a suitable choice of  $p, q$  it is possible to ensure that  $\|w^\dagger - \mathbb{P}_p w^\dagger\|_{\ell^2} \leq c_p \delta$  and  $\|(I - \mathbb{P}_q)A\|_{X \rightarrow Y} \leq c_q \delta$ . Then, provided that  $\lambda$  is chosen as  $\lambda = c_0 \delta$ , there exists a positive constant  $c_2$  (depending on  $\|A\|, \|w^\dagger\|_{\ell^1}$ , on the choice of  $\{\psi_I\}, \{\varphi_j\}$ , and on the constants  $c_0, c_1, c_p, c_q$ ) such that*

$$(2.6) \quad \|w_{\delta,p,q} - w^\dagger\|_{\ell^1} \leq c_2 \delta.$$

The proof, which follows by an application of the variational source condition reported in [17, section 3], is reported in Appendix A.

**Remark 2.3.** Upper bounds of the kind  $\|w^\dagger - \mathbb{P}_p w^\dagger\|_{\ell^2} \leq f(p)$  can be explicitly computed under some particular assumptions on  $w^\dagger$ . If, for example, we suppose that  $u^\dagger$  is a cartoon-like image (i.e.,  $u^\dagger$  is a  $C^2$ -smooth function apart from a jump discontinuity along a finite set of  $C^2$ -curves) and  $\{\psi_I\}$  is the Haar wavelets basis, it is well known that  $\|w^\dagger - \mathbb{P}_p w^\dagger\|_{\ell^2} \leq p^{-1}$  (see, e.g., [36, Chapter 9]).

On the other hand, an estimate for the term  $\|(I - \mathbb{P}_q)A\|_{X \rightarrow Y}$  can be obtained by standard results of finite-rank approximation of compact operators. For example, suppose that the operator  $A$  is a compact operator. Define  $\{s_j\}$  as its singular values (i.e., let  $\{(s_j, e_j)\}$  be the eigenvalues and eigenfunctions of  $(A^*A)^{\frac{1}{2}}$ ) and suppose the sequence  $s_j$  is nonincreasingly converging to 0. A sufficient condition for this is that  $A$  is a Schatten operator of any class  $p$ . If we select the basis  $\{\varphi_j\}$  such that  $\varphi_j = U e_j$ , where  $U$  is the partial isometry in the polar decomposition  $A = U(A^*A)^{\frac{1}{2}}$ , then it holds that  $\|(I - \mathbb{P}_q)A\|_{X \rightarrow Y} \leq s_{q+1}$ . In the case of the Radon transform in 2 dimensions, according to [38, section IV.3],  $s_j = c_R j^{-\frac{1}{2}}$ , hence to get  $\|(I - \mathbb{P}_q)A\|_{X \rightarrow Y} \leq c_R \delta$  it is enough to consider  $q \geq \frac{1}{\delta^2} - 1$ .

A well-know technique for the solution of the minimization problem (2.5) is the ISTA (introduced in [12]), which consists in selecting an initial guess  $w^{(0)} \in \mathbb{R}^p (\cong W_p)$  and in iteratively creating the sequence  $\{w^{(n)}\}_{n=1}^N$  as follows:

$$(2.7) \quad w^{(n)} = \mathcal{T}(w^{(n-1)}) = \mathcal{S}_{\lambda/L} \left( w^{(n-1)} - \frac{1}{L} W A_{p,q}^* A_{p,q} W^* w^{(n-1)} + \frac{1}{L} W A_{p,q}^* m \right),$$

where  $\frac{1}{L} > 0$  is interpreted as a (fictitious) time step and, for  $\beta > 0$ ,  $S_\beta(w)$  is the (component-wise) soft-thresholding operator:

$$[S_\beta(w)]_I = S_\beta(w_I); \quad S_\beta(w_I) = \begin{cases} w_I + \beta & \text{if } w_I < -\beta, \\ 0 & \text{if } |w_I| \leq \beta, \\ w_I - \beta & \text{if } w_I > \beta. \end{cases}$$

The convergence of  $\{w^{(N)}\}$  to a minimizer  $w_{\delta,p,q}$  of (2.5) is analyzed, in an infinite-dimensional context, in [8]. The following result for the discrete problem under consideration is instead a direct consequence of [6, Theorem 25].

**Proposition 2.4.** *If  $L$  is chosen such that  $L \geq \|WA_{p,q}^*A_{p,q}W^*\|/2$  then the sequence  $\{w^{(N)}\}$  generated via (2.7) by any  $w^{(0)} \in \mathbb{R}^p$  converges in  $\ell^2$  to the solution  $w_{\delta,p,q}$  of (2.5). Moreover, there exist  $c_3 > 0$  and  $0 \leq a < 1$  (both depending on  $A_{p,q}$ ,  $L$ , and  $\|w^\dagger\|_{\ell^2}$ ) such that*

$$(2.8) \quad \|w^{(N)} - w_{\delta,p,q}\|_{\ell^2} \leq c_3 a^N.$$

**2.2. A modification of ISTA.** We now consider a perturbation of ISTA (2.7). Let  $Z : \ell^2(\mathbb{N}) \rightarrow \ell^2(\mathbb{N})$  satisfy

$$(2.9) \quad \|WA_{p,q}^*A_{p,q}W^* - Z\|_{\ell^2 \rightarrow \ell^2} \leq \rho.$$

Then, we substitute  $Z$  in place of the matrix  $WA_{p,q}^*A_{p,q}W^*$  in the expression of ISTA. To note the dependency on the perturbation amplitude  $\rho$ , we denote by  $\{w_\rho^{(n)}\}$  the sequence obtained by selecting  $w_\rho^{(0)} \in \mathbb{R}^p$  and iterating

$$(2.10) \quad w_\rho^{(n)} = \mathcal{T}_Z(w_\rho^{(n-1)}) = \mathcal{S}_{\lambda/L} \left( w_\rho^{(n-1)} - \frac{1}{L} Z w_\rho^{(n-1)} + \frac{1}{L} W A_{p,q}^* m \right).$$

The following result shows a connection between the convergence of the sequence  $\{w_\rho^{(n)}\}$  to the minimizer  $w_{\delta,p,q}$  and the magnitude of the perturbation  $\rho$ .

**Proposition 2.5.** *Let  $w^{(0)} = w_\rho^{(0)}$ ,  $L \geq \|WA_{p,q}^*A_{p,q}W^*\|$ , and consider  $N_0, \eta_0 > 0$ . Then there exists a constant  $\tilde{c}_4$ , depending on  $L, A, w^{(0)}, \|w^\dagger\|_{\ell^2}$ , and on  $N_0, \eta_0$ , such that if  $N \geq N_0$  and  $\rho N \leq \eta_0$ , then*

$$(2.11) \quad \|w_\rho^{(N)} - w_{\delta,p,q}\|_{\ell^2} \leq c_3 a^N + \tilde{c}_4 \rho N.$$

*If, moreover,  $N, \rho$  are chosen as  $N > \frac{\ln(\delta^{-1})}{\ln(a^{-1})}$  and  $\rho < \frac{\delta}{N}$ , then (for  $c_4 = c_3 + \tilde{c}_4$ )*

$$(2.12) \quad \|w_\rho^{(N)} - w_{\delta,p,q}\|_{\ell^2} \leq c_4 \delta.$$

The proof of this proposition follows by the nonexpansivity of the soft-thresholding operator and is reported in [Appendix B](#).

We collect the results obtained in [Propositions 2.2](#) and [2.5](#) in the following final convergence estimate.

**Theorem 2.6.** *Let  $w^\dagger$  satisfy (2.1) and let  $A$  be injective. For sufficiently small  $\delta$ , select a regularization parameter  $\lambda = c_0\delta$ . Select  $p, q$  such that  $\|w^\dagger - \mathbb{P}_p w^\dagger\| \leq c_p\delta$  and  $\|(I - \mathbb{P}_q)A\|_{X \rightarrow Y} \leq c_q\delta$ . Let  $L \geq \|WA_{p,q}^* A_{p,q} W^*\|$  and consider the perturbed ISTA iterations (2.10), where the operator  $Z$  satisfies (2.9),  $N = \log_a \delta$ , and  $\rho = \frac{\delta}{N}$ . Then, there exists a positive constant  $c_5$  (depending on the previously introduced constants  $c_0, c_1, c_2, c_3, c_4, c_p, c_q$ ) such that, for sufficiently small  $\delta$ ,*

$$(2.13) \quad \|w_\rho^{(N)} - w^\dagger\|_{\ell^2} \leq c_5\delta.$$

**2.3. Wavelets in 2 dimensions.** In order to derive the main results of the paper, we need to assume that the orthogonal basis  $\{\psi_I\}_{I=1}^\infty$  is a wavelet basis in  $X = L^2(\Omega)$ . Although our approach is sufficiently general to handle higher-dimensional spaces, we are going to focus on the two-dimensional case, i.e.,  $\Omega \subset \mathbb{R}^2$  (e.g.,  $\Omega = [0, 1]^2$ ). Before moving to the representation of the operator  $A^*A$  with respect to such a basis, we need to describe in more details its structure.

A common way to define a wavelet basis in  $\mathbb{R}^2$  is to rely on two real functions  $\psi$  and  $\phi$ , respectively, defined as mother wavelet and scaling function, whose support is in  $[0, 1]$ . We identify an element  $\psi_I$  of the basis by its scale  $j$ , its translation  $k \in \mathbb{N}_0^2$ , and its type  $(t) \in \{(v), (h), (d), (f)\}$  (respectively, vertical, horizontal, diagonal, and low-pass filter). We denote  $\psi_I(x)$  as  $\psi_{j,k}^{(t)}(x) = 2^j \psi^{(t)}(2^j x - k)$ ,  $x \in [0, 1]^2$ , where we have

$$\begin{aligned} \psi^{(v)}(x_1, x_2) &= \phi(x_1)\psi(x_2), & \psi^{(h)}(x_1, x_2) &= \psi(x_1)\phi(x_2), \\ \psi^{(d)}(x_1, x_2) &= \psi(x_1)\psi(x_2), & \psi^{(f)}(x_1, x_2) &= \phi(x_1)\phi(x_2). \end{aligned}$$

When selecting a maximum scale  $J$  (and  $J_0 < J$  as coarsest scale), we can define a wavelet basis of  $p = 2^{2J}$  elements as follows: take  $j \in \{J_0, \dots, J_1 = J - 1\}$ ; for each  $j \neq J_0$ , consider wavelets of the types  $(v)$ ,  $(h)$ , and  $(d)$ , whereas for  $j = J_0$  include also the type  $(f)$ . For each level  $j$  and type  $(t)$ , consider offsets  $k = (k_1, k_2)$ ,  $k_1 = 0, \dots, 2^j - 1$ ,  $k_2 = 0, \dots, 2^j - 1$ .

We group the wavelet basis functions in subbands, each of which is identified by a scale  $j$  and a type  $(t)$ , obtaining  $3(J - J_0) + 1$  subsets.

**3. ISTA and CNNs.** It is already well known that the unrolled iterations of ISTA can be considered as the layers of a neural network (see, e.g., [22]). Indeed, the  $n$ th iteration of ISTA can be written as

$$(3.1) \quad w^{(n)} = \mathcal{S}_{\lambda/L} \left( w^{(n-1)} - \frac{1}{L} K^{(n)} w^{(n-1)} + \frac{1}{L} b^{(n)} \right),$$

being that  $K^{(n)} = WA_{p,q}^* A_{p,q} W^*$  and  $b^{(n)} = WA_{p,q}^* m$ , independently of  $n$ . At the same time, (3.1) can be seen as the  $n$ th layer of a recurrent neural network, where  $K^{(n)}$  is the matrix of the weight coefficients and  $b^{(n)}$  is the bias vector. Notice that the resulting architecture is the one of a recurrent neural network although, due to its theoretical deduction, it does not present any advanced residual block (such as skip connections) which are common features in the related literature. We also point out that formula (3.1) enforces a specific choice of the nonlinear activation function, namely, the soft-thresholding operator, instead of the more widely used



rectified linear unit (ReLU) or sigmoid functions. Additionally, the soft thresholding operator  $S_\alpha$  can be written in terms of ReLUs as follows:

$$(3.2) \quad S_\alpha(x) = \max(0, x - \alpha) - \max(0, -x - \alpha) = \text{ReLU}(x - \alpha) - \text{ReLU}(-x - \alpha)$$

for  $x \in \mathbb{R}$ ; for vectors, it is applied componentwise.

When considering only the first  $N$  iterations of ISTA, we can collect the parameters appearing in the layers in a vector  $\theta \in \Theta$ . Together with the entries of the matrices  $K^{(n)}$ , we may consider as parameters the step length  $L$  as well as the regularization parameter  $\lambda$ ; see [subsection 5.3](#) for more details. Conversely, the bias vectors  $b^{(n)}$  are not to be considered as parameters: they are fixed and equal to  $WA_{p,q}^*m$  in each layer. We then introduce the map  $f_\theta : Y \rightarrow \ell^1(\mathbb{N})$ , parameterized by  $\theta \in \Theta$ , which takes as an input  $m \in Y_q$  and computes  $N$  iterations like (3.1), where, for each  $n$ ,  $K^{(n)} \in \mathbb{R}^{p \times p}$  is specified in  $\theta$  and  $b^{(n)} = WA_{p,q}^*m$ . For any selected value of  $p, q, N, \lambda, L$ , we know that there exists a particular choice  $\theta_0$  which corresponds to the ISTA iterations associated with the measurements  $m$ .

In this section we show that, under some assumptions on the operator  $A$ , it is possible to interpret the operations in (3.1) as a layer of a CNN. We therefore provide a fairly general network architecture which allows one to recover the standard ISTA iterations (or a perturbation of the kind described by (2.9)), for a specific choice of the parameters.

From now on, we focus on the case  $X = L^2(\Omega)$ , and consider a wavelet basis  $\{\psi_I\}$  of the kind described in [subsection 2.3](#).

**3.1. A convolutional interpretation of ISTA.** We first show, under additional assumptions on operator  $A$ , how to translate the neural network encoded by the operator  $f_\theta$  above into a CNN, allowing for a significant reduction of the number of parameters involved. Suppose that  $A^*A$  is a convolutional kernel operator, i.e.,

$$(3.3) \quad \mathbf{K}_{I,I'} = (A^*A\psi_I, \psi_{I'})_X = \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} K(x, x') \psi_I(x) \psi_{I'}(x') dx dx',$$

$$K(x, x') = K(x - x').$$

According to the description in [subsection 2.3](#), the wavelet basis can be naturally split into subbands, each of which is identified by a couple  $j, (t)$ . This implies that the matrix  $\mathbf{K}$  representing  $A^*A$  can be seen as a block matrix. We now aim at describing the application of each block  $\mathbf{K}_{j \rightarrow j'}^{(t) \rightarrow (t')}$  by means of the following operations:

1. *Discrete convolution.* Let  $B \in \mathbb{R}^{b \times b}$ ,  $C \in \mathbb{R}^{(2b-1) \times (2b-1)}$ , and denote the elements of  $C$  with indices  $i, j$  with  $i = -b + 1, \dots, 0, \dots, b - 1$ ,  $j = -b + 1, \dots, 0, \dots, b - 1$ . Then,  $C * B \in \mathbb{R}^{b \times b}$ :

$$(3.4) \quad (C * B)_{k,l} = \sum_{i=0}^{b-1} \sum_{j=0}^{b-1} C_{k-i, l-j} B_{i,j}.$$

2. *Upsampling.* Let  $B \in \mathbb{R}^{b \times b}$ , then,  $\mathcal{U}(B) \in \mathbb{R}^{2b \times 2b}$  satisfies

$$(3.5) \quad \mathcal{U}(B)[2k : 2k + 1, 2l : 2l + 1] = \begin{bmatrix} B_{k,l} & 0 \\ 0 & 0 \end{bmatrix} \quad \forall k, l = 0, \dots, b - 1,$$

where the notation  $\mathcal{U}(B)[2k : 2k + 1, 2l : 2l + 1]$  is used to denote a submatrix of  $\mathcal{U}(B)$  containing the rows from  $2k$  to  $2k + 1$  and all the columns from  $2l$  to  $2l + 1$ . We denote by  $\mathcal{U}^\eta$  the iterated application of  $\mathcal{U}$ :  $\mathcal{U}^\eta = \mathcal{U} \circ \dots \circ \mathcal{U}$  ( $\eta$  times).

3. *Downsampling.* Let  $B \in \mathbb{R}^{2b \times 2b}$ , then,  $\mathcal{D}(B) \in \mathbb{R}^{b \times b}$  satisfies

$$(3.6) \quad \mathcal{D}(B)_{k,l} = B_{2k,2l} \quad \forall k, l = 0, \dots, b - 1.$$

We denote by  $\mathcal{D}^\eta$  the iterated application of  $\mathcal{D}$ :  $\mathcal{D}^\eta = \mathcal{D} \circ \dots \circ \mathcal{D}$  ( $\eta$  times).

The following crucial result provides a full description of the convolutional interpretation of the matrix representing  $A^*A$  in the wavelet domain. Such a result can be compared to the ones already known in the literature (see, e.g., [11, formula (4.2)]), although the more complicated structure of the wavelet basis entails some significant differences.

**Proposition 3.1.** *Let  $\mathbf{K} \in \mathbb{R}^{p \times p}$  be the matrix representing an operator  $A^*A$  satisfying (3.3) in a two-dimensional (2D) wavelet basis  $\{\psi_I\}_{I=1}^p$ . For a vector  $w \in \mathbb{R}^p$ , let  $w_j^{(t)}$  be the vector of the wavelet components related to basis functions of scale  $j$  and type  $(t)$ . Let  $\mathbf{K}_{j \rightarrow j'}^{(t) \rightarrow (t')}$  denote the block of  $\mathbf{K}$  corresponding to the  $j, (t)$  subset of the column indices and the  $j', (t')$  subset of the row indices. Then*

$$(3.7) \quad \mathbf{K}_{j \rightarrow j'}^{(t) \rightarrow (t')} w_j^{(t)} = \begin{cases} \mathcal{D}^\delta(\tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')} * W_j^{(t)}) & \text{if } j > j', \\ \tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')} * W_j^{(t)} & \text{if } j = j', \\ \tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')} * \mathcal{U}^\delta(W_j^{(t)}) & \text{if } j < j' \end{cases}$$

with  $\delta = |j' - j|$  and  $\tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')} \in \mathbb{R}^{(2^{\hat{j}+1}-1) \times (2^{\hat{j}+1}-1)}$  (where  $\hat{j} = \max(j, j')$ ):

$$(3.8) \quad \left[ \tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')} \right]_d = \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} K(x - x' - 2^{-\hat{j}}d) \psi_{j',0}^{(t')}(x') \psi_{j,0}^{(t)}(x) dx dx', \\ d = (d_1, d_2), \quad d_1, d_2 = \{-2^{\hat{j}} + 1, \dots, 0, \dots, 2^{\hat{j}} - 1\}.$$

The matrix  $W_j^{(t)} \in \mathbb{R}^{2^j \times 2^j}$  is obtained by reshaping the vector  $w_j^{(t)} \in \mathbb{R}^{2^{2j}}$  so that  $[W_j^{(t)}]_d$  is the component  $w_I$  whose index is identified by  $(j, (t), d)$ .

*Proof.* Let  $I, I'$  be identified by  $(j, (t), k)$  and  $(j', (t'), k')$ , respectively. Then,

$$\begin{aligned} [\mathbf{K}]_{I',I} &= \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} K(x - x') \psi_{j,k'}^{(t')}(x') \psi_{j,k}^{(t)}(x) dx dx' \\ &= \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} K(x - x') \psi_{j,0}^{(t')}(x' - 2^{-j'}k') \psi_{j,0}^{(t)}(x - 2^{-j}k) dx dx' \\ &= \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} K(x + 2^{-j}k - x' - 2^{-j'}k') \psi_{j,0}^{(t')}(x) \psi_{j,0}^{(t)}(x) dx dx' \\ &= \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} K(x - x' - 2^{-\hat{j}}(2^{\delta^-}k' - 2^{\delta^+}k)) \psi_{j,0}^{(t')}(x) \psi_{j,0}^{(t)}(x) dx dx' = \left[ \tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')} \right]_d, \end{aligned}$$

where  $\delta^+ = \max(0, j - j')$ ,  $\delta^- = \max(0, j' - j)$ , and  $d = 2^{\delta^-}k' - 2^{\delta^+}k$ . For the sake of ease, we use  $\mathbf{K}$  instead of  $\mathbf{K}_{j \rightarrow j'}^{(t) \rightarrow (t')}$ ,  $\tilde{\mathbf{K}}$  instead of  $\tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')}$ ,  $w$  instead of  $w_j^{(t)}$ ,  $W$  instead of  $W_j^{(t)}$ .

Moreover, we denote by  $\mathcal{I}$  the set of indices  $\mathcal{I} \subset \{1, \dots, p\}$  belonging to the wavelet scale  $j$  and type  $(t)$ .

Consider first the case  $j = j'$ . Then  $\delta = \delta^+ = \delta^- = 0$ , and it holds

$$[\mathbf{K}]_{I',I} = [\tilde{\mathbf{K}}]_d, \quad d = k' - k.$$

Therefore,

$$\begin{aligned} [\mathbf{K}w]_{I'} &= \sum_{I \in \mathcal{I}} [\mathbf{K}]_{I',I} w_I = \sum_{I \in \mathcal{I}} [\tilde{\mathbf{K}}]_{k'-k(I)} w_I \\ &= \sum_{k_1=-2^j}^{2^j} \sum_{k_2=-2^j}^{2^j} [\tilde{\mathbf{K}}]_{k'_1-k_1, k'_2-k_2} W_{k_1, k_2} = [\mathbf{K} * W]_{I'}. \end{aligned}$$

Let now  $j < j'$ . Then  $\delta = \delta^+ > 0$ ,  $\delta^- = 0$ , and

$$\begin{aligned} [\mathbf{K}w]_{I'} &= \sum_{I \in \mathcal{I}} [\mathbf{K}]_{I',I} w_I = \sum_{I \in \mathcal{I}} [\tilde{\mathbf{K}}]_{k'-2^{\delta^+}k(I)} w_I \\ &= \sum_{k_1=-2^{j'}}^{2^{j'}} \sum_{k_2=-2^{j'}}^{2^{j'}} [\tilde{\mathbf{K}}]_{k'_1-2^{\delta^+}k_1, k'_2-2^{\delta^+}k_2} \mathcal{U}^{\delta^+}(W)_{2^{\delta^+}k_1, 2^{\delta^+}k_2} = [\mathbf{K} * \mathcal{U}^{\delta^+}W]_{I'}. \end{aligned}$$

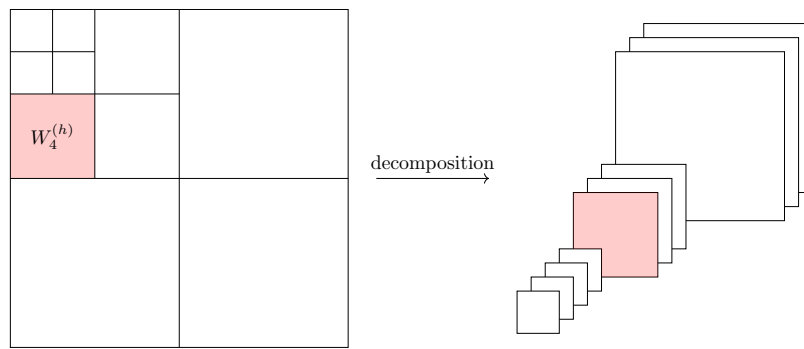
Finally, let  $j > j'$ . Then  $\delta^+ = 0$ ,  $\delta = \delta^- > 0$ , and

$$\begin{aligned} [\mathbf{K}w]_{I'} &= \sum_{I \in \mathcal{I}} [\mathbf{K}]_{I',I} w_I = \sum_{I \in \mathcal{I}} [\tilde{\mathbf{K}}]_{2^{\delta^-}k'-k(I)} w_I \\ &= \sum_{k_1=-2^j}^{2^j} \sum_{k_2=-2^j}^{2^j} [\tilde{\mathbf{K}}]_{2^{\delta^-}k'_1-k_1, 2^{\delta^-}k'_2-k_2} W_{k_1, k_2} = [\mathcal{D}^{\delta^-}(\mathbf{K} * W)]_{I'}. \end{aligned}$$

*Remark 3.2.* The most relevant consequence of [Proposition 3.1](#) is a significant reduction of the number of coefficients required to describe the application of  $A^*A$  as a function from  $\mathbb{R}^p$  to  $\mathbb{R}^p$ . The standard representation, obtained by a matrix in  $\mathbb{R}^{p \times p}$ , indeed involves  $p^2 = 2^{4J}$  parameters, whereas the representation via the convolutional filters  $\tilde{\mathcal{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')}$  involves only  $O(p)$  elements.

This convolutional interpretation also reflects on the neural network architecture proposed in [\(3.1\)](#): if we substitute the multiplication  $K^{(n)}w^{(n-1)}$  by the operations encoded by [\(3.7\)](#) (decomposition of  $w^{(n-1)}$  in wavelet subbands, upscaling, application of convolutional filters, downscaling), the parameters  $\theta$  involved in the description of  $K^{(n)}$  are reduced. The representation of the linear operators  $K^{(n)}$  through convolutions, upscaling, and downscaling is a typical feature of CNNs; thus, by designing a CNN which reproduces exactly the operations reported in [\(3.7\)](#) and [\(3.1\)](#), we can ensure that such a network is completely equivalent, for a suitable choice  $\theta_0$  of the parameters, to the application of ISTA.

**3.2. A working example.** In order to better visualize the convolutional representation of ISTA reported in [\(3.8\)](#), we now provide a small example. Consider the case of  $64 \times 64$  images, thus associated with  $J = 6$  and  $p = 2^{12}$ . Create a wavelet basis consisting of three scales of



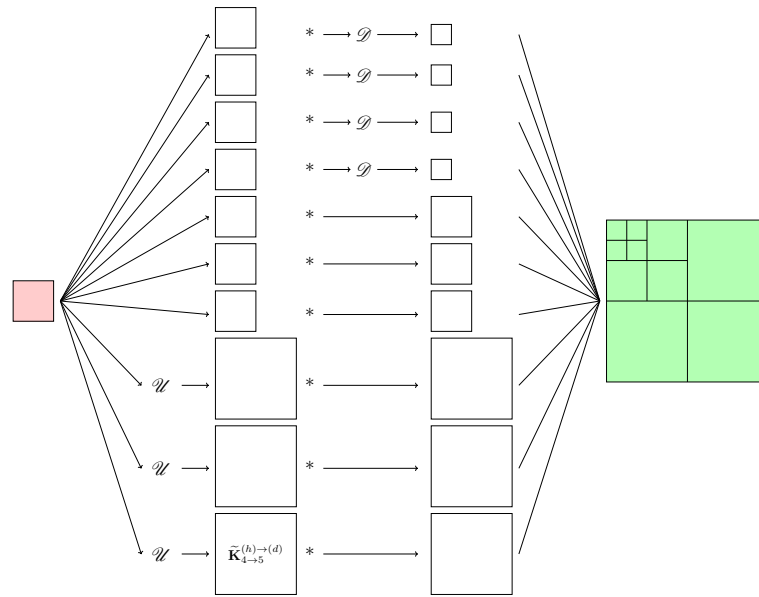
**Figure 1.** Interpretation of (3.7). Step 1: decompose the wavelet transform into subbands.

wavelets, from  $J_0 = 3$  to  $J_1 = 5$ . The resulting basis  $\{\psi_I\}_{I=1}^p$  can therefore be split into 10 subbands: 4 associated with the scale  $j = 3$  (types  $(h)$ ,  $(v)$ ,  $(d)$ , and  $(f)$ ); 3 associated with the scales  $j = 4$  (types  $(h)$ ,  $(v)$ ,  $(d)$ ) and 3 with  $j = 5$ . Each subband consists of  $2^{2j}$  elements. The operator  $A^*A$  is represented in the wavelet basis  $\{\psi_I\}$  by a matrix  $\mathbf{K} \in \mathbb{R}^{p \times p}$ . According to subsection 3.1, the following procedure is equivalent to applying the matrix  $\mathbf{K}$  on a vector  $w \in \mathbb{R}^p$  (representing the wavelet transform of an image):

1. First, split the vector  $w$  into its 10 wavelet subbands, each of which is identified by a scale  $j$  and a type  $(t)$ . This operation is depicted in Figure 1. The vector  $w_j^{(t)} \in \mathbb{R}^{2^{2j}}$  can also be interpreted as a matrix  $W_j^{(t)} \in \mathbb{R}^{j \times j}$ . The element  $[W_j^{(t)}]_d = [W_j^{(t)}]_{(d_1, d_2)}$  is the component associated with the basis function  $\psi_{j,d}^{(t)}(x) = 2^j \psi^{(t)}(2^j x_1 - d_1, 2^j x_2 - d_2)$ .
2. Second, for each subband  $j, (t)$ , compute the 10 vectors  $\mathbf{K}_{j \rightarrow j'}^{(t) \rightarrow (t')} w_j^{(t)}$ , the contributions of  $w_j^{(t)}$  on the subband  $j', (t')$  of the vector  $\mathbf{K}w$ . Each matrix  $\mathbf{K}_{j \rightarrow j'}^{(t) \rightarrow (t')}$  is a  $2^{2j'} \times 2^{2j}$  block composing the matrix  $\mathbf{K}$ . According to (3.7), this can be done by means of usampling, downsampling, and convolution. Consider the case  $j = J_0 = 3$ :
  - if  $j' = 3$ , then  $\hat{j} = 3$  and  $\delta = 0$ . Thus, if we compute the convolution of the  $15 \times 15$  filter  $\tilde{\mathbf{K}}_{3 \rightarrow 3}^{(t) \rightarrow (t')}$  with the matrix  $W_3^{(t)} \in \mathbb{R}^{8 \times 8}$ , we get a  $8 \times 8$  matrix representing the vector  $\mathbf{K}_{3 \rightarrow 3}^{(t) \rightarrow (t')} w_3^{(t)} \in \mathbb{R}^{64}$ .
  - if  $j' = 4$ , then we shall use the third variant in formula (3.7) with  $\delta = 1$  (whereas in (3.8) we have  $\hat{j} = 4$ ). To compute the  $16 \times 16$  matrix associated with  $\mathbf{K}_{3 \rightarrow 4}^{(t) \rightarrow (t')} w_3^{(t)}$ , we must first upsample the matrix  $W_3^{(t)}$  and then convolve it with the  $31 \times 31$  filter  $\tilde{\mathbf{K}}_{3 \rightarrow 4}^{(t) \rightarrow (t')}$ .
  - if  $j' = 5$ , then we again use the third variant of (3.7), with  $\delta = 2$ ; hence the matrix  $W_3^{(t)}$  must be upsampled twice before being convolved with the  $63 \times 63$  filter  $\tilde{\mathbf{K}}_{3 \rightarrow 5}^{(t) \rightarrow (t')}$ .

Consider instead the case  $j = 4$ :

- if  $j' = 3$ , then we need to use the first variant in (3.7) with  $\delta = 1$  (and (3.8) with  $\hat{j} = 4$ ), which means we first compute the convolution between the  $31 \times 31$  filter  $\tilde{\mathbf{K}}_{4 \rightarrow 3}^{(t) \rightarrow (t')}$  and the matrix  $W_4^{(t)} \in \mathbb{R}^{16 \times 16}$  and then downscale it to recover the  $8 \times 8$  matrix describing  $\mathbf{K}_{4 \rightarrow 3}^{(t) \rightarrow (t')} w_4^{(t)}$ .



**Figure 2.** Interpretation of (3.7). Step 2: convolution, upsampling, and downsampling.

- the case  $j' = 4$  is analogous to  $3 \rightarrow 3$ , using  $31 \times 31$  filters  $\tilde{\mathbf{K}}_{4 \rightarrow 4}^{(t) \rightarrow (t')}$ .
- the case  $j' = 5$  is analogous to  $3 \rightarrow 4$ : we first perform upsampling and then convolution.

Finally, for  $j = J_1 = 5$ ,

- if  $j' = 3$ , then we first compute the convolution between  $\tilde{\mathbf{K}}_{5 \rightarrow 3}^{(t) \rightarrow (t')} \in \mathbb{R}^{63 \times 63}$  and  $W_5^{(t)} \in \mathbb{R}^{32 \times 32}$  and then downsample twice.
- if  $j' = 4$ , we only downsample once, as in the case  $4 \rightarrow 3$ .
- if  $j' = 5$ , we only do convolution, as in the cases  $3 \rightarrow 3$  and  $4 \rightarrow 4$ , but with  $63 \times 63$  filters.

A graphical visualization of these operations is provided by Figure 2.

3. The last step consists of collecting, for each subband  $j', (t')$ , all the contributions coming from the vectors  $w_j^{(t)}$ . Thanks to the previous step, among the 100 computed matrices, all the ones associated with those contributions have dimensions  $2^{j'} \times 2^{j'}$ . By adding them up we recover the  $j', (t')$  subband of the vector  $\mathbf{K}w$  (see Figure 3).

**3.3. On the possibility of using smaller filters.** When designing a CNN, it is common practice to employ a large numbers of convolutional filters of small size. In the architecture determined by (3.7) and (3.8), the required number of filters is exactly  $(3(J - J_0) + 1)^2$ , and each part of the vector  $w^{(n-1)}$  interacts only with  $(3(J - J_0) + 1)$  of them. Moreover, the size of each filter must be equal to  $(2^{j'+1} - 1)(2^{j'+1} - 1)$ . We now consider the effect of substituting for such large filters with smaller ones.

We would like to use filters of size  $\tau \times \tau$ , with  $\tau = (2\xi + 1)$  and  $\xi > 1$ , obtained by extracting the central elements of the large filters  $\tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')}$ . In particular, we define  $\tilde{\mathbf{K}}^\tau = (\tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')})_\tau$

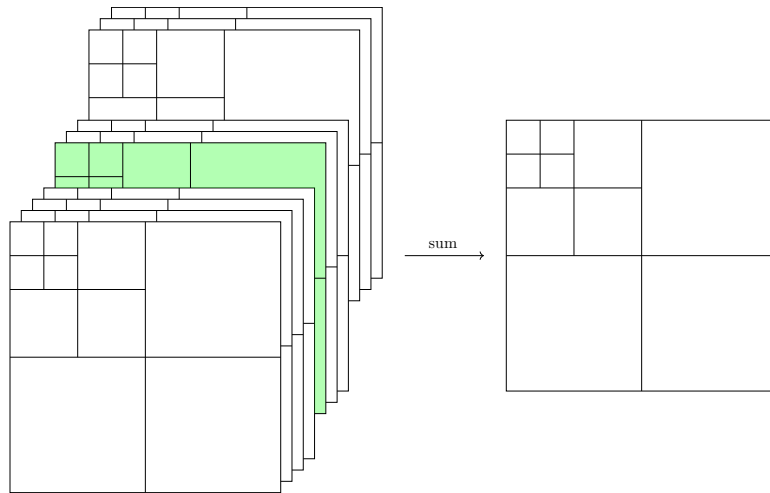


Figure 3. Interpretation of (3.7). Step 3: reassembling each wavelet subband.

with  $\tau = 2\xi + 1$ , as

$$(3.9) \quad [\tilde{\mathbf{K}}^\tau]_d = \begin{cases} [\tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')}]_d & \text{if } \|d\|_\infty \leq \xi, \\ 0 & \text{if } \|d\|_\infty > \xi. \end{cases}$$

We claim that this modification is equivalent to performing a perturbation of ISTA of the type treated in Proposition 2.5, where the parameter  $\rho$  is a suitable function of  $\tau$ . Although providing a detailed proof of this would entail cumbersome computation, we prove the most important result which is required to accomplish this task: we exhibit a bound on the coefficients of the filters which are discarded due to (3.9).

Such an estimate can be obtained by assuming further hypotheses on the operator  $A$ . In particular, suppose that  $A^*A$  is a convolutional operator of kernel  $K$  (as in (3.3)) and, in addition, that for  $x \neq x'$  the kernel  $K(x, x')$  is smooth and such that

$$(3.10) \quad K(x, x') \leq \frac{C}{|x - x'|}, \quad |\nabla_x K(x, x')| + |\nabla_{x'} K(x, x')| \leq \frac{C}{|x - x'|^2}.$$

It is easy to verify that (3.10) is satisfied whenever  $A^*A$  is a  $\Psi$ DO of order  $-1$  with constant coefficients, that is,

$$A^*A f = \mathcal{F}^{-1} \{a(\xi) \mathcal{F} \{f\}(\xi)\}, \quad a(\xi) \sim \frac{1}{|\xi|} \text{ as } \xi \rightarrow 0.$$

We also assume the first-order vanishing moment property for wavelet basis functions:

$$(3.11) \quad \int_{\mathbb{R}^2} \psi_I(x) dx = 0.$$

Such a property is verified even by 2D Haar wavelets, apart from the type ( $f$ ).

**Proposition 3.3.** *Let the operator  $A$  satisfy (3.3) and (3.10). Let the indices  $I, I'$  denote two wavelets of scales  $j, j'$ , type  $(t), (t')$ , and offsets  $k, k'$ . Let  $\psi_I$  and  $\psi_{I'}$  satisfy (3.11) and let  $d_{I, I'}$  be the distance between the supports of  $\psi_I$  and  $\psi_{I'}$ . Whenever  $d_{I, I'} > 0$ , it holds that*

$$(3.12) \quad \mathbf{K}_{I, I'} = (A^* A \psi_I, \psi_{I'})_X \leq c \frac{2^{-2(j+j')}}{d_{I, I'}^3}.$$

We remark that the decay reported in (3.12) closely resembles formula (9.22) in [10] (according to the choice  $n = 2$ ,  $\tilde{d} = 1$ ,  $r = 2t = -1$ ) and with minor changes also formula (4.26) in [5] (with  $M = 2$ ).

*Proof.* According to (3.10), and to (3.11), for any choice of  $x_0 \in \text{supp } \psi_I$ ,  $x'_0 \in \text{supp } \psi_{I'}$  there exists two points  $\xi, \xi'$  in the same supports such that

$$\begin{aligned} \mathbf{K}_{I, I'} &= \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} (K(x, x') - K(x, x'_0)) \psi_I(x) \psi_{I'}(x') dx dx' \\ &\leq \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} |\nabla_x K(x, \xi)| |x' - x'_0| \psi_I(x) \psi_{I'}(x') dx dx' \\ &\leq C \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} \frac{|x' - x'_0|}{|x - \xi|^2} \psi_I(x) \psi_{I'}(x') dx dx' \leq C \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} \frac{|x - x_0| |x' - x'_0|}{|\xi - \xi'|^3} \psi_I(x) \psi_{I'}(x') dx dx'. \end{aligned}$$

The quantity  $|\xi - \xi'|$  is bounded from below by  $d_{I, I'}$  by definition. Moreover,  $|x - x_0| \leq \text{diam}(\text{supp } \psi_I) = c2^{-j}$  and, finally,  $\int_{\mathbb{R}^2} \psi_I(x) \leq 2^j |\text{supp } \psi_I| = 2^{-j}$  (analogous arguments hold for  $I'$ ). ■

In view of (3.12) and of (3.8), we can easily obtain a bound on the elements of the convolutional filters:

$$\left[ \tilde{\mathbf{K}}_{j \rightarrow j'}^{(t) \rightarrow (t')} \right]_d \leq c \frac{2^{-\hat{j}}}{(\|d\|_\infty - 1)^3},$$

provided that  $\|d\|_\infty > 1$ . This result, together with (3.7), allows one to obtain an explicit bound (in the form of (2.12)) on the perturbation induced by the thresholding (3.9).

**4.  $\Psi$ DONet: Formulation and theoretical results.** In this section we introduce a reconstruction algorithm for sparsity-promoting regularization based on CNNs, which leads to a novel network architecture defined as  $\Psi$ DONet. We report the general idea inspiring such a technique, taking advantage of the theoretical results obtained in section 3 and providing a comprehensive interpretation. Eventually, we provide a theoretical result ensuring the convergence of the proposed algorithm.

**4.1.  $\Psi$ DONet: A network to learn  $\Psi$ DOs.** Inspired by the results of section 3, if the operator  $A^*A$  is of convolutional type, we define a reconstruction algorithm by designing a CNN of  $N$  layers, each of which is described by (3.1). In particular, the bias vectors appearing in (3.1) are  $b^{(n)} = WA_{p,q}^* m$  for each  $n$ , whereas the linear operators  $K^{(n)}$  are interpreted as a combination of upscaling, downscaling, and convolution as described in (3.7). As shown in Proposition 3.1, if the entries of the convolutional filters are selected as is (3.8), this procedure is equivalent to performing  $N$  iterations of ISTA. Instead, the key idea of

the proposed algorithm is to split the convolutional filters into two parts: a central  $\tau \times \tau$  submatrix (where  $\tau$  is a predefined hyperparameter of the algorithm) and the outer frame. For each one of the  $3(J - J_0) + 1$  filters required by each layer, we suppose that the entries in the external frame are specified according to (3.8), whereas the central entries are considered as parameters, to be learned throughout the training process. Such parameters are collected in a vector  $\theta_n$  (related to the  $n$ th layer) and ultimately stored in the vector  $\theta$ , possibly together with other learnable parameters. The obtained network is denoted as  $f_\theta^\tau$ : the aim of a CNN-based algorithm is to find a parameter  $\theta$  such that the network is a good approximation of the solution map of our inverse problem, taking as an input the measurements  $m$  and giving as an output the solution  $w^\dagger = Wu^\dagger$ .

It is evident that, among the possible choices of the optimal parameter, the network could select the vector  $\theta_0$  which exactly replicates the ISTA iterations (it is the one for which, in every layer, the central entries of each filter are also specified by (3.8)). Nonetheless, if the optimal choice of  $\theta$  differs from  $\theta_0$ , it means that the network is learning something more than the ISTA iterations associated with the operator  $A^*A$ . This can be meaningfully interpreted as follows: in each layer, the network  $f_\theta^\tau$  applies the filters associated with an operator whose kernel is  $K_0 + K_1$ , where  $K_0$  is the kernel of  $A^*A$  and  $K_1$  is the kernel of another, *learned*, operator. Since the difference will only occur in the central elements of the convolutional filters, according to the analysis of subsection 3.3, we can argue that the learned operator is indeed a suitable approximation of a  $\Psi$ DO. This finally allows us to motivate the name we propose for this novel CNN-based reconstruction algorithm:  $\Psi$ DONet.

There are several reasons for which the learning process could attain a better result than the one provided by ISTA. Indeed, a better choice of the parameters allows us to reduce numerical errors induced by the discrete representation of  $A^*A$ , which might have a significant effect due to the error propagation among the iterations. Moreover, we might also mitigate model errors in the definition of the operator  $A$  itself. Finally, this perturbation could provide a representation of  $A^*A$  with respect to a slightly different basis, which allows us to better satisfy the sparsity assumption on the solutions. For such reasons, the use of  $\Psi$ DONet is specifically recommended whenever the original operator  $A^*A$  is a  $\Psi$ DO itself. Indeed, its kernel representation by means of convolutional filters might benefit from learned corrections in all its most important entries, namely, the central ones.

We will show that  $\Psi$ DONet is also highly recommended for FIOs: in this case, the largest entries of the convolutional filters representing  $A^*A$  are located in the center and along some lines, possibly stretching away from the center. This is the case of the limited-angle Radon transform (deeply analyzed in the following sections), which is associated with the kernel

$$K(x, y) = \frac{1}{|x - y|} \chi_\Gamma(x - y),$$

with  $\chi_\Gamma$  the indicator function of the cone in  $\mathbb{R}^2$  between the angles  $-\Gamma$  and  $\Gamma$ . As reported in section 5, the convolutional filters related to this operator show large values only in the central elements and along two lines having the same slope as the ones delimiting the cone. This provides a curious shape for the filters, which resemble a *bow tie*. We will show that the application of  $\Psi$ DONet on this operator, providing learned corrections only to the center of the bow ties, is still extremely effective.



In addition to the numerical verification of the previous statements (depicted in section 6), we provide here a theoretical argument to explain why  $\Psi$ DONet is expected to outperform ISTA even when  $A^*A$  is a FIO.

**4.2. A theoretical justification for  $\Psi$ DONet using microlocal analysis.** If a  $\Psi$ DO correction is learned for the normal operator, the modified ISTA iterations read as

$$w^{(n)} = \mathcal{S}_{\lambda/L} \left( w^{(n-1)} - \frac{1}{L} W(A_{p,q}^* A_{p,q} + \Lambda_p) W^* w^{(n-1)} + \frac{1}{L} W A_{p,q}^* m \right),$$

with  $\Lambda_p$  the discrete representation of a  $\Psi$ DO. This is equivalent to performing the standard ISTA iterations on a modified version of a minimization problem (2.5), namely,

$$\min_{w \in \ell^1(\mathbb{N})} \|(A_{p,q}^* A_{p,q} + \Lambda_p)^{1/2} W^* w - (A_{p,q}^* A_{p,q} + \Lambda_p)^{-1/2} A_{p,q}^* m\|_X^2 + \lambda \|w\|_{\ell^1}.$$

This minimization problem is a discretised version of the continuous minimization problem

$$\min_{u \in Z} \|(A^*A + \Lambda)^{1/2} u - (A^*A + \Lambda)^{-1/2} A^* m\|_Y^2 + \lambda \|u\|_Z.$$

Eventually, this amounts to finding a regularized solution, with a regularization penalty promoting solutions for which  $w = Wu \in \ell^1(\mathbb{N})$  is sparse, of the problem

$$(A^*A + \Lambda)^{1/2} u = (A^*A + \Lambda)^{-1/2} A^* m$$

or, equivalently,

$$(4.1) \quad (A^*A + \Lambda)u = A^* m.$$

Assume next that  $A^*A$  is a FIO that defines a bounded map between Sobolev spaces in a ball  $B(R)$  of radius  $R$  and that there is  $r \in \mathbb{R}$  such that for all  $s \in \mathbb{R}$ ,  $A^*A : H_0^s(B(R)) \rightarrow H^{s+r}(B(R))$ . Let  $\Lambda$  be a (possibly unbounded) self-adjoint, positive definite, and invertible operator  $\Lambda : L^2(B(R)) \rightarrow L^2(B(R))$ . Moreover, assume that  $\Lambda$  is given by an elliptic  $\Psi$ DO of order  $h$  and  $h > r$ . Then the operator  $A^*A\Lambda^{-1}$  is an operator smoothing Sobolev spaces by order  $r + h$ , that is,  $A^*A\Lambda^{-1} : H^s(B(R)) \rightarrow H^{s+r+h}(B(R))$  for all  $s \in \mathbb{R}$ . Moreover, the operator  $(A^*A + \Lambda)^{-1}$  can be written as

$$(A^*A + \Lambda)^{-1} = \Lambda^{-1} (I + A^*A\Lambda^{-1})^{-1} = \Lambda^{-1} \left( \sum_{k=0}^{\bar{k}} (-A^*A\Lambda^{-1})^k + R_{\bar{k}} \right),$$

where  $R_{\bar{k}} : H^s \rightarrow H^{s+\bar{k}(r+h)}$  is bounded for all  $s \in \mathbb{R}$ . Thus the solution  $u_\Lambda$  of (4.1) can be written as

$$u_\Lambda = \sum_{k=0}^{\bar{k}} u_k + \Lambda^{-1} R_{\bar{k}} A^* m, \quad u_k = -\Lambda^{-1} (A^*A\Lambda^{-1})^k = (-\Lambda^{-1} A^*A)^k \Lambda^{-1} A^* m.$$

Note that here the operator  $\Lambda^{-1} A^*A$  is an FIO whose canonical relation is determined by  $A$  and whose symbol is determined by both  $A$  and  $\Lambda$ . The training of the neural network can

be considered as optimizing  $\Lambda$  so that for a given datum  $m$  the solution  $u_\Lambda$  of (4.1) is close to  $u^\dagger$ . Roughly speaking, this means to optimize  $\Lambda$  so that the imaging artifacts in terms of  $u_k$ , caused by iteration of the operator  $\Lambda^{-1}A^*A$ , are minimized. Note that the remainder term  $\Lambda^{-1}R_{\bar{k}}A^*m$  becomes smoother when  $\bar{k}$  grows.

We are interested in applying this argument in limited-angle tomography. As a first step, we start by considering the case when  $X$ -rays are measured only from finitely many directions  $\omega_j \in S^1$ ,  $j = 1, \dots, J$ . As can be seen, the normal operator obtained via backprojection is associated with the kernel

$$(4.2) \quad K(x, y) = \sum_{j=1}^J \delta((x - y) \cdot \omega_j^\perp),$$

where  $\delta \in \mathcal{D}'(\mathbb{R})$  is the Dirac delta distribution. In this case,  $K = A^*A \in I^\mu(C'_K)$  is an FIO of order  $\mu = -\frac{1}{2}$  and its canonical relation  $C_K = \bigcup_{j=1}^J C_j$ , where

$$C_j = \{(x, \xi; y, \eta) \in (T^*\mathbb{R}^2 \setminus 0) \times (T^*\mathbb{R}^2 \setminus 0) \mid \xi = \eta, \eta \perp \omega_j, (x - y) \cdot \omega_j = 0\},$$

and  $K : H_0^s(B(\mathbb{R})) \rightarrow H^{s+1/2}(B(\mathbb{R}))$  for all  $s \in \mathbb{R}$ . This entails that if  $(y, \eta)$  is in the wavefront set of  $u$ , then the elements  $(x, \xi)$  in the wavefront set of  $Ku$  satisfy  $(x, \xi; y, \eta) \in C_K$ , that is, the operator  $K$  moves singularities along  $C_K$ . This provides a theoretical justification for the appearance of the well-known streaking artifacts in sparse tomography; see [30, 43]. The presence of a  $\Psi$ DO  $\Lambda$  might affect the symbol of the operator, but not the canonical relation of the operator  $\Lambda^{-1}A^*A$ . This amounts to saying that the strength of the singularities of  $u_k = (-\Lambda^{-1}A^*A)^k \Lambda^{-1}A^*m$  can be reduced by  $\Lambda$ , without transporting them.

We finally consider the limited-angle tomography problem, in which the the operator  $K$  is defined as in (1.2). It is possible to show that such an operator belongs to a class of generalized FIO whose properties are studied, e.g., in [21, 20]. In addition,  $K$  can be treated as an FIO if considered far from the diagonal, namely, if a smooth truncation function  $\phi$  is introduced such that  $\phi \in C_0^\infty(\mathbb{R}^2)$  and  $\phi$  vanishes near zero, the kernel  $\phi(x - y)K(x, y)$ , where  $K$  is given in (1.2), defines an FIO. We can thus apply the strategy proposed in subsection 4.1 by supposing to split the kernel  $K = K_0 + K_1$ , where now  $K_0$  is the kernel of  $A^*A$  away from the diagonal and  $K_1$  is the kernel of  $\Lambda$ , a learned correction concentrated on the diagonal. Analogously as above, we see that the operator  $\Lambda$  changes the strength of the artifacts appearing in  $u_k = (-\Lambda^{-1}A^*A)^k \Lambda^{-1}A^*m$  but not their locations. Therefore, even in the limited-angle problem, the  $\Psi$ DO correction can be seen as a regularization technique, from a microlocal analysis standpoint.

**4.3. A convergence result.** We now provide a theoretical result which holds true for the  $\Psi$ DONet algorithm, regardless of its specific implementation. In analogy with the approach of [14], we introduce the following probabilistic approach. Let  $\mathcal{B} = \{u \in X_p : Wu \in \ell^1(\mathbb{N}); \|Wu\|_{\ell^1} \leq C_{\mathcal{B}}\}$  and  $u$  be a random variable having a probability distribution  $\mu$  on the space  $\mathcal{B}$ . We can consider  $\mu$  as some prior information on the solution of the inverse problem. Moreover, let  $\epsilon$  be a random variable in  $Y_q$  with distribution  $\nu$ , which models the error on the measurements. Assume that  $u$  and  $\epsilon$  are independent, hence, the measurement  $m = A_{p,q}u + \epsilon$  is a random variable on the product space  $X_p \times Y_q$  with density  $A_*\mu \otimes \nu$ , where  $A_*\mu$  denotes

the pushforward of  $\mu$  to  $Y$  via the linear map  $A$ . In order to measure the performance of the network  $f_\theta^\tau$ , we introduce the loss function associated with the network  $f_\theta^\tau$  as

$$(4.3) \quad \mathcal{L}(\theta; \mu, \nu) = \mathbb{E}_{u \sim \mu, \epsilon \sim \nu} [\|f_\theta^\tau(A_{p,q}u + \epsilon) - Wu\|_{\ell^2}^2].$$

We define the optimal neural network as the one associated with  $\theta^*$  satisfying

$$(4.4) \quad \theta^* = \arg \min_{\theta \in \Theta} \mathcal{L}(\theta; \mu, \nu).$$

Before focusing on the properties of the optimal network  $f_{\theta^*}^\tau$ , it is convenient to recall that, for a specific choice of parameters  $\theta_0$ , the network  $f_{\theta_0}^\tau$  is equivalent to performing  $N$  iterations of (modified) ISTA. The following rough estimate will be useful.

**Lemma 4.1.** *There exist two constants  $k_1, k_2 > 0$  (depending on  $C_B, L, \rho, \|A_{p,q}\|, w^{(0)}, N$ ) such that, for all  $u \in \mathcal{B}$  and  $\epsilon \in Y_q$ ,*

$$(4.5) \quad \|f_{\theta_0}^\tau(A_{p,q}u + \epsilon) - Wu\|_{\ell^2} \leq k_1 + k_2\|\epsilon\|_{Y_q}.$$

*Proof.* According to (3.1), defining  $\kappa = 1 + \frac{\|A_{p,q}^* A_{p,q}\| + \rho}{L}$ , we get

$$\begin{aligned} \|f_{\theta_0}^\tau(A_{p,q}u + \epsilon) - Wu\|_{\ell^2} &\leq \|f_{\theta_0}^\tau(A_{p,q}u + \epsilon)\|_{\ell^2} + \|u\|_{X_p} \\ &\leq \kappa^N \|w^{(0)}\| + (1 + \kappa + \dots + \kappa^{N-1}) \|A_{p,q}u + \epsilon\|_{Y_q} + C_B \\ &\leq \kappa^N \|w^{(0)}\| + C_B + \frac{\kappa^N - 1}{\kappa - 1} (\|A_{p,q}\| C_B + \|\epsilon\|_{Y_q}). \quad \blacksquare \end{aligned}$$

We now focus on the case in which  $\epsilon$  is a Gaussian random vector, i.e.,  $\nu = N(0, \sigma^2 I_q)$ , with  $I_q$  the identity matrix in  $\mathbb{R}^{q \times q}$ . In this case, it is useful to recall that

$$(4.6) \quad \mathbb{E}[\|\epsilon\|_{Y_q}^2] = q\sigma^2, \quad \mathbb{E}[\|\epsilon\|_{Y_q}^4] \leq 3q^2\sigma^4.$$

In addition to Lemma 4.1, we can rely on the results reported in section 2 (and in particular on Theorem 2.6) to provide a more refined estimate. Indeed, we observe that the convergence result reported in (2.13) is independent of the choice of  $\epsilon = m - Au^\dagger$ , as long as  $\|\epsilon\| \leq \delta$ . Moreover, the constant  $c_5$  appearing in (2.13) can depend on  $u^\dagger$ , but only through an upper bound on  $\|w^\dagger\|_{\ell^1}$  (see, in particular, [17, Theorem 1] and [6, Theorem 25] for the constant derived from Proposition 2.2 and Proposition 2.5, respectively). This allows us to conclude the following.

**Lemma 4.2.** *Suppose  $\epsilon \sim N(0, \sigma^2 I_q)$  and let  $\delta = \sigma^{1/\eta}$ , with  $\eta > 1$ . There exists  $\sigma_0 > 0$  such that, for  $\sigma < \sigma_0$ , for every  $u \in \mathcal{B}$*

$$\mathbb{E}_{\epsilon \sim \nu} [\|f_{\theta_0}^\tau(Au + \epsilon) - Wu\|_{\ell^2}^2] \leq c_5^2 \delta^2 + 2\sqrt{2}k_1^2 \delta^{\eta-1} + 2\sqrt{6}k_2^2 q \delta^{3\eta-1}.$$

*If, moreover,  $\eta = 3$  and  $\sigma < \min\{\sigma_0, q^{-1/2}\}$ , then there exists a constant  $c^*$  (depending on  $c_5, k_1, k_2$ ) such that*

$$(4.7) \quad \mathbb{E}_{\epsilon \sim \nu} [\|f_{\theta_0}^\tau(Au + \epsilon) - Wu\|_{\ell^2}^2] \leq c^* \delta^2.$$

*Proof.* We start by considering that

$$\begin{aligned} \mathbb{E}_{\epsilon \sim \nu} [\|f_{\theta_0}^\tau(Au + \epsilon) - Wu\|_{\ell^2}^2] &= \int_{Y_q} \|f_{\theta_0}^\tau(Au + \epsilon) - Wu\|_{\ell^2}^2 d\nu(\epsilon) \\ &= \int_{\|\epsilon\| < \delta} \|f_{\theta_0}^\tau(Au + \epsilon) - Wu\|_{\ell^2}^2 d\nu(\epsilon) + \int_{\|\epsilon\| > \delta} \|f_{\theta_0}^\tau(Au + \epsilon) - Wu\|_{\ell^2}^2 d\nu(\epsilon). \end{aligned}$$

We now employ (2.13) on the first term and the Hölder inequality on the second term. Moreover, in view of Chebyshev’s inequality,  $\nu(\{\|\epsilon\| > \delta\}) \leq \frac{\sigma^2}{\delta^2}$ . Therefore,

$$\begin{aligned} \mathbb{E}_{\epsilon \sim \nu} [\|f_{\theta_0}^\tau(Au + \epsilon) - Wu\|_{\ell^2}^2] &\leq c_5^2 \delta^2 \left(1 - \frac{\sigma^2}{\delta^2}\right) + \frac{\sigma}{\delta} \left(\mathbb{E}_{\epsilon \sim \nu} [\|f_{\theta_0}^\tau(Au + \epsilon) - Wu\|_{\ell^2}^4]\right)^{\frac{1}{2}} \\ &\leq c_5^2 \delta^2 + \frac{\sigma}{\delta} \left(8k_1^4 + 8k_2^4 \mathbb{E}[\|\epsilon\|_{Y_q}^4]\right)^{\frac{1}{2}}. \end{aligned}$$

By (4.6) and by  $\sigma = \delta^\eta$  we immediately verify the first thesis, and imposing  $\eta = 3$  and  $\delta^{2\eta} q < 1$  we get (4.7) with  $c^* = c_5^2 + 2\sqrt{2}k_1^2 + 2\sqrt{6}k_2^2$ . ■

In view of Lemma 4.2, we can easily prove the following convergence result regarding the optimal network  $f_{\theta^*}^\tau$ .

**Proposition 4.3.** Consider  $\epsilon \sim N(0, \sigma^2 I_q)$  with  $\delta = \sigma^{1/3}$  and let  $\theta^*$  satisfy (4.4). There exists  $\sigma_1 > 0$  such that, for  $\sigma \leq \min\{\sigma_1, q^{1/2}\}$ , it holds that

$$(4.8) \quad \mathcal{L}(\theta^*; \mu, \nu) \leq c^* \delta^2.$$

This also amounts to saying that the random variable  $f_{\theta^*}^\tau(A_{p,q}u + \epsilon)$  converges to  $Wu$  in the mean as  $\delta \rightarrow 0$ .

*Proof.* By the definition of  $\theta^*$  and by Lemma 4.2,

$$\begin{aligned} \mathcal{L}(\theta^*; \mu, \nu) &\leq \mathcal{L}(\theta_0; \mu, \nu) = \int_{\mathcal{B}} \int_{Y_q} \|f_{\theta_0}^\tau(A_{p,q}u + \epsilon) - Wu\|_{\ell^2}^2 d\nu(\epsilon) d\mu(u) \\ &\leq \int_{\mathcal{B}} c^* \delta^2 d\mu(u) = c^* \delta^2. \end{aligned}$$

Although the optimal network  $f_{\theta^*}^\tau$  allows for a precise approximation of the solution map of the inverse problems, it is impractical for solving the minimization problem stated in (4.4). Instead, neural network algorithms require one to draw a sample from the random variables  $U$  and  $E$  and to find the parameter  $\theta$  which allows for the best reconstruction on such a sample (training process). This task is addressed by minimizing a discretized loss functional, as reported in subsection 5.4, and results in the definition of the *trained* neural network. The quality of the trained network can be verified by analyzing its *generalization*, namely, its ability to provide good predictions even when tested on data outside the training sample. Such an analysis has been performed in detail (although with some different assumptions with respect to the ones in this work) in [14], and can be extended also to the problem under consideration.

**5. In practice: The particular case of CT.** In this section, we focus on the practical aspects of the reconstruction algorithm introduced in [subsection 4.3](#), in the particular case of LA-CT with the discrete setting. In the remainder of the article, the discrete counterpart of the operator  $A_{p,q}$  representing the LA-CT will be denoted by  $R_\Gamma$ . We first define the regularized minimization problem, and then propose an effective method for the computation of the convolutional kernel filters approximating the backprojection operator in the wavelet domain. Third, we present and discuss the general reconstruction workflow and finally give more details on the two CNN architectures we propose in this paper.

**5.1. The CT minimization problem.** After suitable discretization, we are given the measurements (i.e., the so-called sinogram or observed image)  $\mathbf{m} \in \mathbb{R}^q$  such that

$$(5.1) \quad \mathbf{m} = \mathbf{R}_\Gamma \mathbf{u}^\dagger + \epsilon,$$

where  $\mathbf{u}^\dagger \in \mathbb{R}^p$  denotes the (unknown) discrete and vectorized image,  $\mathbf{R}_\Gamma \in \mathbb{R}^{q \times p}$  describes a discretized version of the Radon transform where the angles are limited in the arc specified by  $\Gamma$ , and  $\epsilon \in \mathbb{R}^q$  models the measurement noise. We call  $\mathbf{w}^\dagger$  the  $\mathbb{R}^p$ -vector such that  $\mathbf{W}\mathbf{u}^\dagger = \mathbf{w}^\dagger$ , where  $\mathbf{W} \in \mathbb{R}^{p \times p}$  represents a discretization of the wavelet transform. Thus, the regularized minimization problem is given by

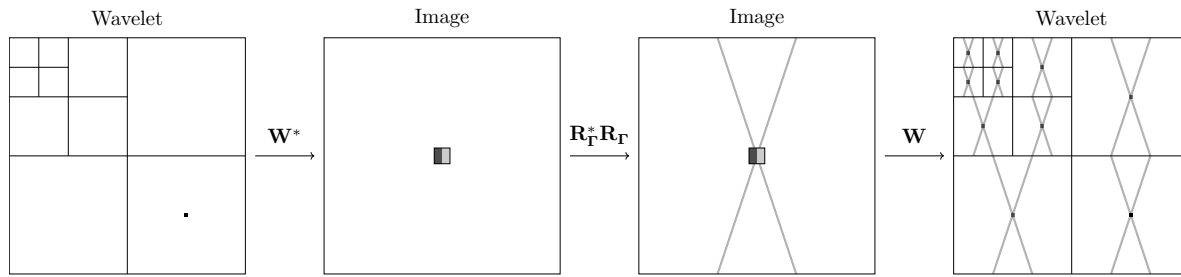
$$(5.2) \quad \min_{\mathbf{w} \in \mathbb{R}^p} \|\mathbf{R}_\Gamma \mathbf{W}^* \mathbf{w} - \mathbf{m}\|_2^2 + \lambda \|\mathbf{w}\|_1.$$

Our recovery algorithm for finding a reconstruction  $\mathbf{u}$  of  $\mathbf{u}^\dagger$  involves convolutional architectures incorporated into the iterative structure of standard ISTA, as described in [section 3](#). In the next paragraphs, we detail the implementation of such an algorithm.

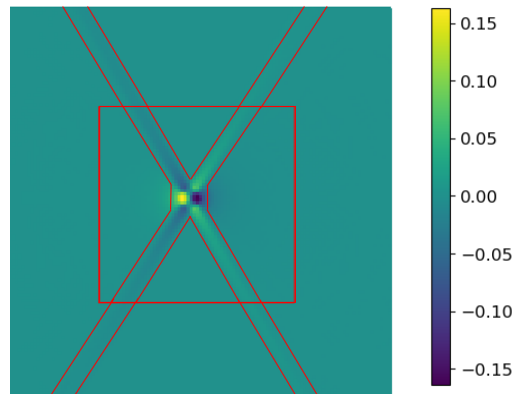
**5.2. Convolutional kernel operator for limited-angle CT.** Building a convolutional algorithm that reproduces the behavior of standard ISTA first requires identifying the various blocks of the matrix  $\mathbf{K}$  representing the backprojection operator in the wavelet domain. In other words, the very first step in the development of our method is to establish the convolutional filters of  $\mathbf{K}$  which, once applied as defined in [\(3.7\)](#), provide a reliable approximation of the operator  $\mathbf{W}\mathbf{R}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$ .

One way to compute such convolutional filters that proves to be a numerically advantageous alternative to [\(3.8\)](#), is represented in [Figure 4](#). Let us consider an object whose representation in the wavelet domain has only one nonzero pixel, located at the center of one of its wavelet subbands. Applying the operator  $\mathbf{W}\mathbf{R}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$  to this initial object leads to a new object whose subbands present a bow-tie-shaped structure. Those “bow-tie” subbands constitute a first set of convolutional filters. By reiterating this operation until the central pixel of all the wavelet subbands in the initial object has been visited, one obtains the entire collection of convolutional filters necessary for the approximation of  $\mathbf{W}\mathbf{R}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$ . A numerical example of convolutional filter is shown in [Figure 5](#).

In order to imitate the behavior of the operator  $\mathbf{W}\mathbf{R}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$ , the convolutional filters so computed are then to be applied to the wavelet subbands of an object as illustrated in [Figure 6](#). First, each wavelet subband of the object of interest is replicated  $3(J - J_0) + 1$  times. Those

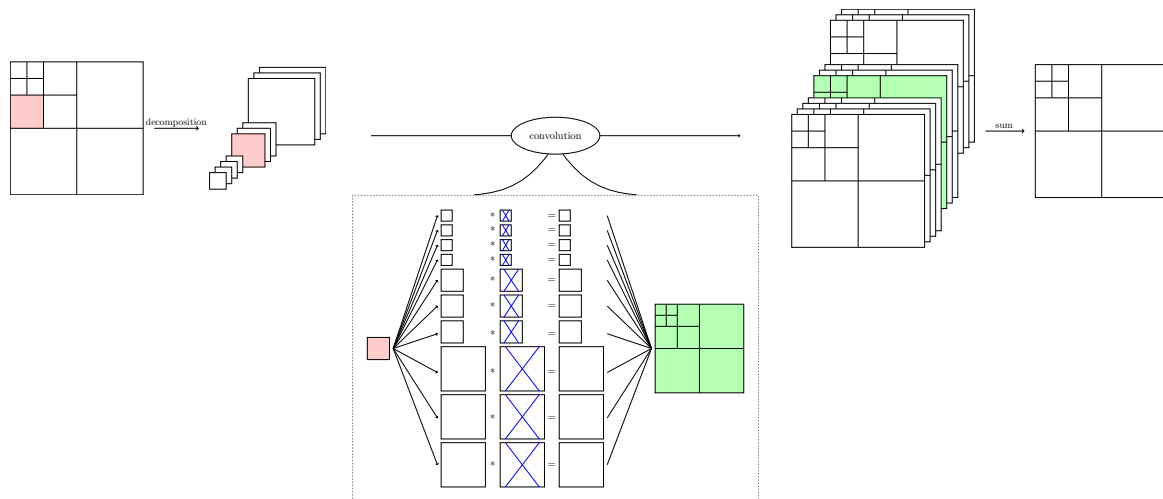


**Figure 4.** Illustration of the proposed way to compute the filters of the convolutional kernel operator  $K$  in the LA-CT case. The initial object (on the left) is created such that all its pixels but one are set to zero. The only nonzero pixel is located at the center of one of its wavelet subbands. By applying the operator  $\mathbf{WR}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$  to this initial object, one obtains a new object in the wavelet domain, whose subbands present a bow-tie-shaped structure. The set of bow-tie subbands thus computed from all the possible initial objects constitute the filters of the convolutional kernel operator. Here we have represented three levels of decomposition in the wavelet domain, meaning that the total number of convolutional filters amounts to  $(3^2 + 1)^2 = 100$ .



**Figure 5.** Example of a bow-tie subband that can be used as a convolutional filter of the kernel operator  $K$ . It was generated from a  $256 \times 256$  initial object, according to the procedure detailed in subsection 5.2 and illustrated in Figure 4. Theory suggests that the pixels with highest intensities are spread according to a bow-tie-shaped structure. In practice, they are even more condensed: most of the energy is concentrated along two diagonal lines that intersect in the center and whose inclination is defined by the limited angle: 95.8% of the  $\ell_2$ -norm of the filter is concentrated along those two lines, from which 94.8% are inside the central red square.

replicas are either upsampled, or downsampled, or kept with the same dimensions, depending on the scale of the filter they are to be convolved with. The set of convolutional filters used on the replicas of a particular wavelet subband of scale  $j$  and type  $(t)$  is the set of filters generated beforehand by applying the operator  $\mathbf{WR}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$  to an object whose only non-zero pixel is located at the center of its wavelet subband of the exact same scale  $j$  and type  $(t)$ . Once the convolutions between the replicas and the filters have been performed, the resulting subbands are reassembled to form the wavelet representation of a new object. This process is reiterated for all the subbands of the initial object and, ultimately, the  $3(J - J_0) + 1$  resulting items are summed. The final outcome is an approximation of the wavelet representation of the operator  $\mathbf{WR}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$  applied to the initial object (cf. Figure 7).



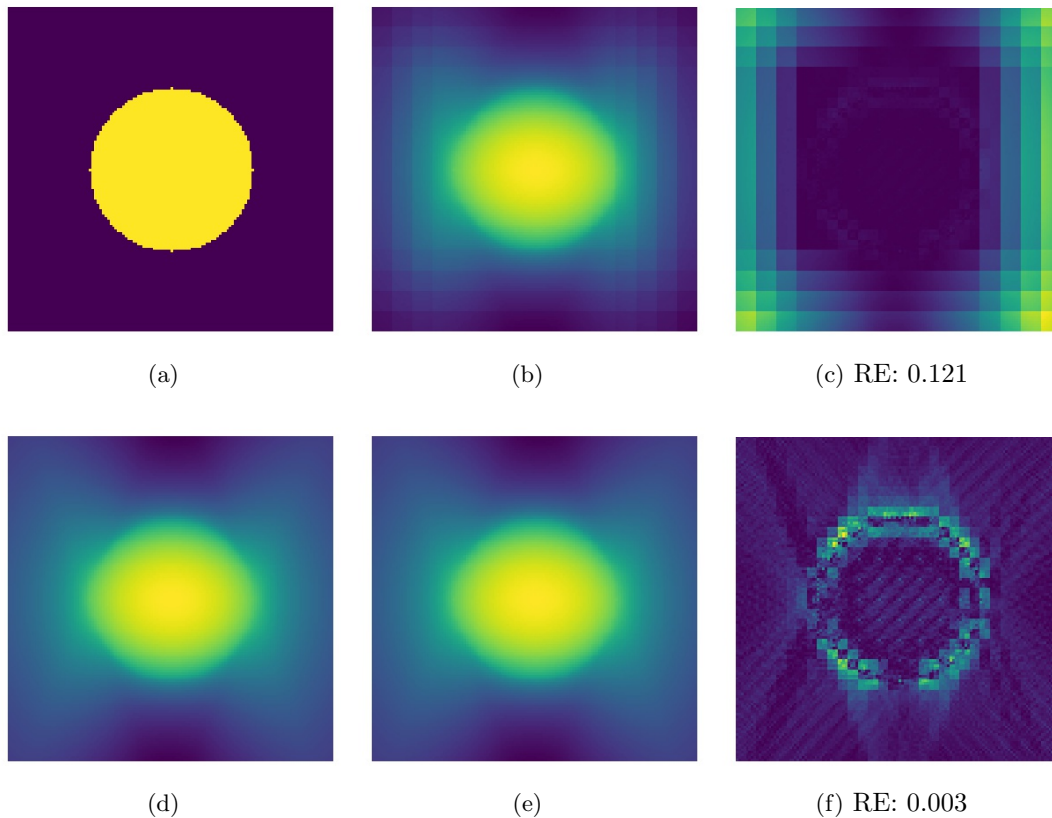
**Figure 6.** Illustration of the way the filters of the convolutional kernel operator are applied to each wavelet subband of the initial object, after up- and down-sampling operations, in order to approximate the operator  $\mathbf{WR}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$ .

Two remarks are worth mentioning regarding the creation and use of the above-defined convolutional filters. First, our practical implementation very slightly differs from the theory presented in (3.7) as far as the downsampling is concerned. In our codes, downsampling is indeed applied before computing the convolution between the filter and the wavelet subband replica, and not after as it is presented in the theory. This choice is motivated by the reduction in terms of storage needs and running time such a change allows while preserving the accuracy of the approximation. Second, both the theoretical analysis and the experimental tests showed that the dimensions of the convolutional filters used for the approximation of  $\mathbf{WR}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$  do affect the accuracy of the results. Initially, we assumed that the convolutional filters should be generated from only-one-nonzero-pixel objects with the same dimensions  $2^J \times 2^J$  as the image of interest (recall that  $p = 2^{2J}$ ). However, we reached the conclusion that they actually have to be generated from twice bigger objects, that is of dimensions  $2^{J+1} \times 2^{J+1}$ , in order to get an accurate approximation of the operator  $\mathbf{WR}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$ . An illustration of the effects of the size of the filters can be seen on Figure 7.

**5.3. Our CNN architectures.** The above described method for generating and applying the filters of the kernel operator  $\mathbf{K}$  makes the concrete implementation of a convolutional algorithm that imitates the behavior of standard ISTA possible. Thus, the convolutional implementation of ISTA, resultwise equivalent to the standard one, could be written as

$$(5.3) \quad \mathbf{w}^{(n+1)} = \mathcal{S}_{\frac{\lambda}{L}} \left( \mathbf{w}^{(n)} + \frac{1}{L} \left( \mathbf{WR}_\Gamma^* \mathbf{m} - \mathbf{Kw}^{(n)} \right) \right), \quad n = \{0, \dots, N\}.$$

This algorithm offers the merits of the iterative model-based method ISTA, while allowing the incorporation of data-driven approaches, such as machine learning and deep neural network techniques. The implementation of  $\mathbf{K}$  indeed involves operations that are all perfectly adaptable to the framework of CNNs. Our goal is precisely to take full advantage of this com-

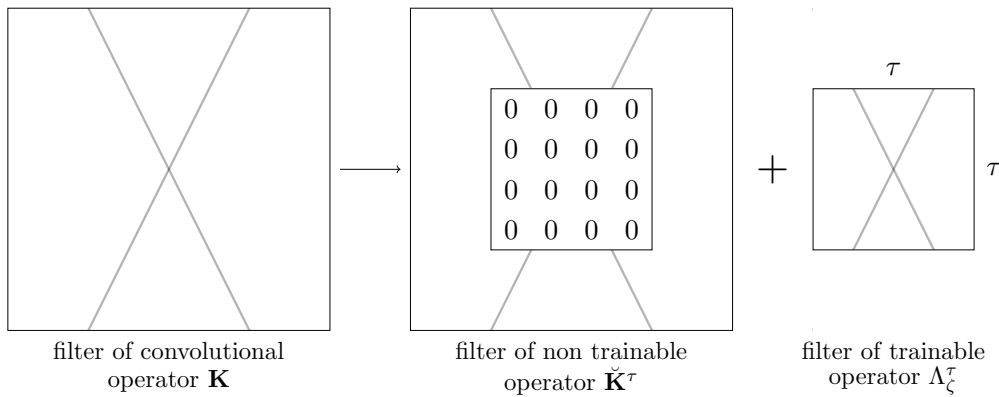


**Figure 7.** Illustration of the effect of the standard backprojection operator and of its approximations based on convolutional filters of different sizes. (a) shows the ground truth  $\mathbf{u}^\dagger$  of interest and (d) its standard back-projection  $\mathbf{R}_F^\dagger \mathbf{R}_F \mathbf{u}^\dagger$ , computed with the basic functions of the Python package *scikit-image*. (b) (resp., (e)) represents the approximation  $\mathbf{K} \mathbf{u}^\dagger$  obtained with convolutional filters beforehand generated from  $2^J \times 2^J$  (resp.,  $2^{J+1} \times 2^{J+1}$ ) only-one-nonzero-pixel object. (c) and (f) show the absolute differences between the approximation of the backprojection operator and the expected value (d). The dynamic range of the plot is modified for better contrast.

patibility and profit from the remarkable potentials of deep neural networks by converting the hitherto fixed operator  $\mathbf{K}$  into a partially trainable CNN. Thus, the center of the convolutional filters so far precomputed with the deterministic method presented in subsection 5.2 can henceforth be considered as parameters to be learned from data. The choice of learning only the central part of the convolutional filters of  $\mathbf{K}$  rather than the whole filters is motivated by the need to reduce the model complexity which, in the latter case, makes the training of the model burdensome if not impractical.

In order to further improve reconstruction performance, we also propose to learn the soft-thresholding parameter as well as the step length so far set at  $1/L$ . The so-defined convolutional architecture results in our proposed algorithm  $\Psi$ DONet, whose convergence results are detailed in subsection 4.3. In subsections 5.3.1 and 5.3.2, we propose two different implementations of  $\Psi$ DONet, that prove to be resultwise similar as can be observed in subsection 6.2.





**Figure 8.** Illustration of the way the convolutional filters of the two operators in  $\Psi\text{DONet-O}$  are computed based on the filters of operator  $\mathbf{K}$ . Each filter of  $\mathbf{K}$  is partitioned into two filters, which sum is equivalent to the initial one. The filter of  $\check{\mathbf{K}}^\tau$  is a copy of the filter of  $\mathbf{K}$  with the exception that the  $\tau \times \tau$  central weights are set to zero. The filter of  $\Lambda_\zeta^\tau$  has dimensions  $\tau \times \tau$  and is initialized with the central  $\tau \times \tau$  central weights of the filter of  $\mathbf{K}$ .

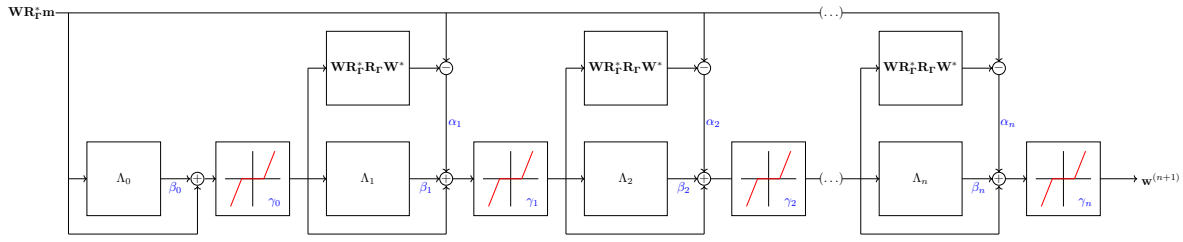
**5.3.1.  $\Psi\text{DONet-F}$ .** The most natural way to implement  $\Psi\text{DONet}$  consists in partitioning the convolutional operator  $\mathbf{K}$  into two operators: a fixed one,  $\check{\mathbf{K}}^\tau$ , and a trainable (single-layer) CNN referred to as  $\Lambda_\zeta^\tau$ , where  $\tau$  is a tunable hyperparameter. The two operators have the exact same architecture as  $\mathbf{K}$  and their sum, before any training, is strictly equivalent to  $\mathbf{K}$ . The first operator  $\check{\mathbf{K}}^\tau$  is nontrainable and its filters are a copy of the filters of  $\mathbf{K}$  with the exception that the  $\tau \times \tau$  central weights of each filter are set to zero. The second operator  $\Lambda_\zeta^\tau$ , on the contrary, is composed by  $\tau \times \tau$ -trainable filters that are initialized with the  $\tau \times \tau$  central part of the filters of  $\mathbf{K}$  (cf. Figure 8). This first implementation of  $\Psi\text{DONet}$ , referred to as filter-based  $\Psi\text{DONet}$  or  $\Psi\text{DONet-F}$ , is formulated as

$$(5.4) \quad \mathbf{w}^{(n+1)} = \mathcal{S}_{\gamma_n} \left( \mathbf{w}^{(n)} + \alpha_n \left( \mathbf{W}\mathbf{R}_\Gamma^* \mathbf{m} - \beta_n \left( \check{\mathbf{K}}^\tau \mathbf{w}^{(n)} + \Lambda_{\zeta_n}^\tau \mathbf{w}^{(n)} \right) \right) \right),$$

where  $n = \{0, \dots, N\}$  and the parameters to be learned are  $\{\gamma_0, \alpha_0, \beta_0, \zeta_0, \dots, \gamma_N, \alpha_N, \beta_N, \zeta_N\}$ . The parameters  $\{\beta_0, \dots, \beta_N\}$  have been added in such a way that the influence of the fixed operator  $\check{\mathbf{K}}^\tau$  with respect to the constant term  $\mathbf{W}\mathbf{R}_\Gamma^* \mathbf{m}$  can be adjusted in order to maximize the accuracy of the results. It is worth mentioning that for the particular choice of  $\gamma_n = \frac{\lambda}{L}, \alpha_n = \frac{1}{L}, \beta_n = 1$  for any  $n = \{0, \dots, N\}$ , this model before any training is exactly equivalent to standard ISTA.

The trade-off between the number of parameters that can be improved through the learning process and the trainability of the model is controlled by  $\tau$ . For a sound choice of such a hyperparameter, the complexity of the model is sufficiently reduced to allow for the convergence of the learning algorithm while enabling the enhancement of a significant number of weights in the filters.

This implementation has the advantage of offering a clear interpretation of the role and meaning of the convolutional filters belonging to  $\check{\mathbf{K}}^\tau$  and  $\Lambda_\zeta^\tau$ . Those filters are indeed initialized with the filters of the operator  $\mathbf{K}$  that imitates the behavior of  $\mathbf{W}\mathbf{R}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$ . Thus, modifying



**Figure 9.** Block diagram of the proposed model (5.5). Notice that the soft-thresholding operator acts as a nonlinear activation function.

their weights through the learning process can be thought of as a direct improvement of the backprojection operator.

$\Psi$ DONet-F has led to very satisfactory preliminary results, presented in section 6. However, training such a model on big images may quickly become extremely onerous in terms of running time and storage requirements. Such problems may arise while training  $\Psi$ DONet-F on images of dimensions greater than or equal to  $256 \times 256$ . Unlike typical CNNs that usually make use of small-sized convolutional filters, the filters of  $\check{\mathbf{K}}^\tau$  in our proposed algorithm are much bigger than the wavelet subbands they are convolved with. This uncommon procedure, that inter alia implies the padding, i.e., the addition of many extra pixels to the edge of each wavelet subband, brings about a severe speed reduction in the training process as well as the necessity of a substantial memory space. The alternative implementation of  $\Psi$ DONet, described in subsection 5.3.2, addresses these shortcomings.

**5.3.2.  $\Psi$ DONet-O.** The main flaw of  $\Psi$ DONet-F rests upon the use of operator  $\check{\mathbf{K}}^\tau$  which implies numerous burdensome convolutions. This issue is worked around in  $\Psi$ DONet-O (5.5), as  $\check{\mathbf{K}}^\tau$  is not involved anymore. Here, the backprojection operator is not approximated, meaning that  $\mathbf{WR}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^*$  is indeed implemented as the succession of the inverse wavelet, Radon, inverse Radon, and direct wavelet transforms applied to the iterate  $\mathbf{w}^{(n)}$ . This second implementation of  $\Psi$ DONet, named operator-based  $\Psi$ DONet or  $\Psi$ DONet-O, reads as

$$(5.5) \quad \mathbf{w}^{(n+1)} = \mathcal{S}_{\gamma_n} \left( \mathbf{w}^{(n)} + \alpha_n \left( \mathbf{WR}_\Gamma^* \mathbf{m} - \mathbf{WR}_\Gamma^* \mathbf{R}_\Gamma \mathbf{W}^* \mathbf{w}^{(n)} \right) + \beta_n \Lambda_{\zeta_n}^\tau \mathbf{w}^{(n)} \right),$$

where  $n = \{0, \dots, N\}$ , the parameters to be learned are  $\{\gamma_0, \alpha_0, \beta_0, \zeta_0, \dots, \gamma_N, \alpha_N, \beta_N, \zeta_N\}$ , and  $\Lambda_{\zeta_n}^\tau$  has the same architecture as the operator  $\mathbf{K}$ . The block diagram of the method is represented in Figure 9. For the special choice of  $\gamma_n = \frac{\lambda}{L}, \alpha_n = \frac{1}{L}, \beta_n = 0$  for any  $n = \{0, \dots, N\}$ , this model is exactly equivalent to standard ISTA. The only convolutions involved in this alternative implementation are the ones composing the CNN  $\Lambda_{\zeta_n}^\tau$ , whose filters are chosen to be small enough to avoid any running time or storage issue. In that sense,  $\Psi$ DONet-O offers an implementation numerically preferable to  $\Psi$ DONet-F, while retaining the same properties on a theoretical level. Furthermore, such a model keeps offering a clear interpretation of its postprocessing abilities since  $\Lambda_{\zeta_n}^\tau$ , on account of its architecture, can still be seen as an adjunct for improving the backprojection operator.

**5.3.3. Note on soft-thresholding parameters.** From a theoretical point of view, the parameters  $\gamma_0, \dots, \gamma_N$  in  $\Psi$ DONet-F and  $\Psi$ DONet-O have to be nonnegative, as they represent the soft-thresholding parameters. In order to stick to the operator originally involved in standard ISTA, it is possible to enforce the positivity of the coefficient by replacing each  $\gamma_n$  by  $\mathbf{10}^{\tilde{\gamma}_n}$ , where  $\tilde{\gamma}_n$  becomes the actual trainable parameter. However, in order to allow for a greater degree of freedom in the learning process, we decided to implement the operator  $\mathcal{S}_{\gamma_n}$  in such a way that it is also interpretable for negative values of its parameter  $\gamma_n$ . In such a case, we define the operator  $\mathcal{S}_{\gamma_n < 0}$  as the symmetric of the soft-thresholding curve with respect to  $\mathbf{y} = \mathbf{x}$ , while for nonnegative values of  $\gamma_n$ ,  $\mathcal{S}_{\gamma_n}$  is exactly equivalent to the soft-thresholding operator. Formally,  $\mathcal{S}_{\gamma_n}$  becomes

$$\mathcal{S}_{\gamma_n}(\mathbf{x}) = \begin{cases} \mathbf{x} - \gamma_n & \text{if } \mathbf{x} \geq \gamma_n, \\ 0 & \text{if } |\mathbf{x}| < \gamma_n, \\ \mathbf{x} + \gamma_n & \text{if } \mathbf{x} \leq -\gamma_n, \end{cases} \quad \text{For } \gamma_n \geq 0 :$$

$$\mathcal{S}_{\gamma_n}(\mathbf{x}) = \begin{cases} \mathbf{x} - \gamma_n & \text{if } \mathbf{x} \geq 0, \\ \mathbf{x} + \gamma_n & \text{if } \mathbf{x} < 0. \end{cases} \quad \text{For } \gamma_n < 0 :$$

The two implementations  $\Psi$ DONet-F and  $\Psi$ DONet-O are tested with and without the positivity constraint on  $\gamma$  (cf. results in [subsection 6.2](#)).

**5.4. Supervised learning.** If we denote  $f_{\theta}^{\tau}$  as the  $N$ -layer CNN that, given  $\mathbf{m}$ , computes the final output  $\mathbf{w}^{(N+1)}$  according to one of the two proposed architectures, we aim at learning the optimal high-dimensional vector  $\theta = \{\gamma_0, \alpha_0, \beta_0, \zeta_0, \dots, \gamma_N, \alpha_N, \beta_N, \zeta_N\}$  that ideally satisfies the relation

$$(5.6) \quad f_{\theta}^{\tau}(\mathbf{m}) \approx \mathbf{W}\mathbf{u}^{\dagger}.$$

For a mathematical formalization, we regard the tuple  $(\mathbf{m}, \mathbf{u}^{\dagger}) \in \mathbb{R}^q \times \mathbb{R}^p$  as a random variable with a joint probability distribution  $\Xi$ , as detailed in [subsection 4.3](#). Ideally, we would like to find a parameter vector  $\theta^*$  minimizing the expected risk

$$(5.7) \quad \min_{\theta} \left( \mathbb{E}_{(\mathbf{m}, \mathbf{u}^{\dagger}) \sim \Xi} \|f_{\theta}^{\tau}(\mathbf{m}) - \mathbf{W}\mathbf{u}^{\dagger}\|_2^2 \right).$$

Other loss functions, such as the weighted  $\mathbf{l}_2$ -norm, where the wavelet coefficients are weighted depending on their scale, have been tested and lead to results similar to the nonweighted  $\mathbf{l}_2$ -norm. For the sake of brevity, we will stick to the basic form of (5.7).

In practice, computing the expectation with respect to  $\Xi$  is not possible. Instead, we are typically given a finite set of independent drawings  $(\mathbf{m}_1, \mathbf{u}_1^{\dagger}), \dots, (\mathbf{m}_S, \mathbf{u}_S^{\dagger})$  and consider the minimization of the empirical risk:

$$(5.8) \quad \min_{\theta} \frac{1}{S} \sum_{i=1}^S \|f_{\theta}^{\tau}(\mathbf{m}_i) - \mathbf{W}\mathbf{u}_i^{\dagger}\|_2^2.$$

Depending on the properties of  $f_{\theta}^{\tau}$ , the optimization problem is in general nonconvex. In the case of neural networks, typically some form of gradient descent is used, where the gradients are calculated via backpropagation [48]. Computing the gradient over the entire training set in (5.8) is often not feasible for large-scale problems due to memory limitations. To circumvent this problem, a stochastic or the minibatch gradient descent is used, in which the gradient is approximated over smaller, randomly selected batches of training examples [19, Chapter 8]. The final performance (i.e., the generalization) of the trained map  $f_{\theta}^{\tau}$  is evaluated on a separate set of independent drawings, the *test set*, that were not previously used in the optimization of  $\theta$  in (5.8).

**6. Experiments and results.** In this section, we evaluate the performance of the proposed reconstruction schemes by comparing the performance with standard ISTA.

**6.1. Preliminaries.** Let us begin by describing the considered experimental scenario, the implementation of the used operators, and the training procedure.

**6.1.1. Data set.** The data set consists of 10700 synthetic images of ellipses, where the number, locations, sizes, and intensity gradients of the ellipses are chosen randomly. Using the MATLAB function `radon`, we simulate measurements for a missing wedge of  $60^{\circ}$  with Gaussian noise. To avoid inverse crime [37] the measurements are simulated at a higher resolution and then downsampled for an image resolution of  $128 \times 128$ . 10000 images are used for training, 200 images for validation, and 500 for testing.

**6.1.2. Operators.** For the implementation of the discrete limited angle operator  $\mathbf{R}_{\Gamma}$  we use the `radon` routine of the Python package `scikit-image` [52], or the 2D parallel beam geometry of the operator discretization library [1], which is based on the Astra toolbox [51]. The former is employed for generating the backprojections  $\mathbf{W}\mathbf{R}_{\Gamma}^* \mathbf{m}$  provided as inputs to  $\Psi$ DONet-F and  $\Psi$ DONet-O, while the latter is used for the implementation of  $\mathbf{W}\mathbf{R}_{\Gamma}^* \mathbf{R}_{\Gamma} \mathbf{W}^*$  in  $\Psi$ DONet-O. The direct and inverse Radon transform operators are multiplied by a constant so that their norm is equal to one. Regarding the wavelet transform, we make use of the Python package `pywt` [34] or a rectified version of the package `tf-wavelets` [23]. In all our experiments, we consider the case  $J = 7$  and  $J_0 = 4$ , implying that the wavelet decomposition  $\mathbf{W}\mathbf{u}$  has 10 subbands. For  $\Psi$ DONet-F and  $\Psi$ DONet-O, we choose to set  $\tau$  to 32. Note that according to theory,  $\tau$  is supposed to be odd, however, in practice we prefer it to be even. This very slight modification has no effect on the results.

**6.1.3. Network structure and training.** For the implementation of  $\Psi$ DONet-F and  $\Psi$ DONet-O, we fix the number of unrolled blocks  $\mathbf{N}$  to 120. In order to reduce the number of parameters to be learned, we choose to use only 40 different sets of trainable parameters  $\{\zeta_n, \gamma_n, \alpha_n, \beta_n\}$ , each of which is being used over 3 consecutive blocks, instead of the theoretically expected 120 sets. Implementing and training our algorithms has been performed using Tensorflow with an Adam optimizer [31] and a learning rate (step size) of  $10^{-3}$ . The number of epochs was chosen to be 3, and the batch size was set to 25. The training, run on a NVIDIA Quadro P6000 GPU, takes roughly 20 hours.<sup>1</sup>

<sup>1</sup>Our codes are available at <https://github.com/megalinier/PsiDONet>.

Table 1

Comparison of reconstruction methods. The similarity values are averaged over the 500 images of the test set.

Method	RE	PSNR	SSIM	HaarPSI
$\mathbf{u}_{\text{ista}}$	0.44	22.84	0.36	0.37
$\mathbf{u}_{\text{FBP}}$	0.64	19.49	0.20	0.30
$\mathbf{u}_{\Psi_{\text{do-F}}}^+$	0.29	26.63	0.59	0.47
$\mathbf{u}_{\Psi_{\text{do-F}}}$	0.25	27.63	0.78	0.54
$\mathbf{u}_{\Psi_{\text{do-O}}}^+$	0.28	26.76	0.60	0.48
$\mathbf{u}_{\Psi_{\text{do-O}}}$	<b>0.23</b>	<b>28.43</b>	<b>0.81</b>	<b>0.58</b>

**6.1.4. Compared methods.** We compare the preliminary results of the architectures we propose with the reconstructions provided by standard ISTA. In the implementation of the latter, we make use of the formula introduced in [12]. The regularization parameter  $\lambda$  and the constant  $L$  are respectively set to  $2 \cdot 10^{-6}$  and 5. The number of iterations for ISTA is determined by the stopping criterion

$$(6.1) \quad \|\mathbf{u}^{(n+1)} - \mathbf{u}^{(n)}\|_2^2 / \|\mathbf{u}^{(n)}\|_2^2 < \text{tol},$$

where  $\text{tol}$  is chosen to be  $2 \cdot 10^{-4}$ . Below, we give a list of the abbreviations henceforth used for the different recovery methods:

- $\mathbf{u}_{\text{ista}}$  Standard ISTA reconstruction.
- $\mathbf{u}_{\text{FBP}}$  Standard FBP with the “ramp” filter of `skicit-image`.
- $\mathbf{u}_{\Psi_{\text{do-F}}}^+$  Solution provided by  $\Psi$ DONet-F with positivity constraint on the soft-thresholding parameter ( $\gamma_n = 10^{\tilde{\gamma}_n} \forall n$ ).
- $\mathbf{u}_{\Psi_{\text{do-F}}}$  Solution provided by  $\Psi$ DONet-F without positivity constraint on the soft-thresholding parameter.
- $\mathbf{u}_{\Psi_{\text{do-O}}}^+$  Solution provided by  $\Psi$ DONet-O with positivity constraint on the soft-thresholding parameter ( $\gamma_n = 10^{\tilde{\gamma}_n} \forall n$ ).
- $\mathbf{u}_{\Psi_{\text{do-O}}}$  Solution provided by  $\Psi$ DONet-O without positivity constraint on the soft-thresholding parameter.

**6.1.5. Similarity measures.** For an assessment of image quality, we are using several quantitative measures, such as the relative error (RE) given by  $\|\mathbf{u}^\dagger - \mathbf{u}\|_2 / \|\mathbf{u}^\dagger\|_2$ , where  $\mathbf{u}^\dagger$  denotes the reference image and  $\mathbf{u}$  its reconstruction. Furthermore, we consider the peak signal-to-noise ratio (PSNR) and the structured similarity index (SSIM) [53] provided by Tensorflow. Finally, we are reporting the Haar wavelet-based perceptual similarity index (HaarPSI) that was recently proposed in [47].

**6.2. Results.** In the following, we will report and discuss the results of our numerical experiments. The average image quality measures of the 500 test images are reported in Table 1. Furthermore, a visualization of the reconstruction quality for two of the test images is given in Figures 10 and 11. Due to the large missing angle of  $60^\circ$ , the FBP images in Figures 10(c) and 11(c) are contaminated with streaking artifacts and contrast changes. The

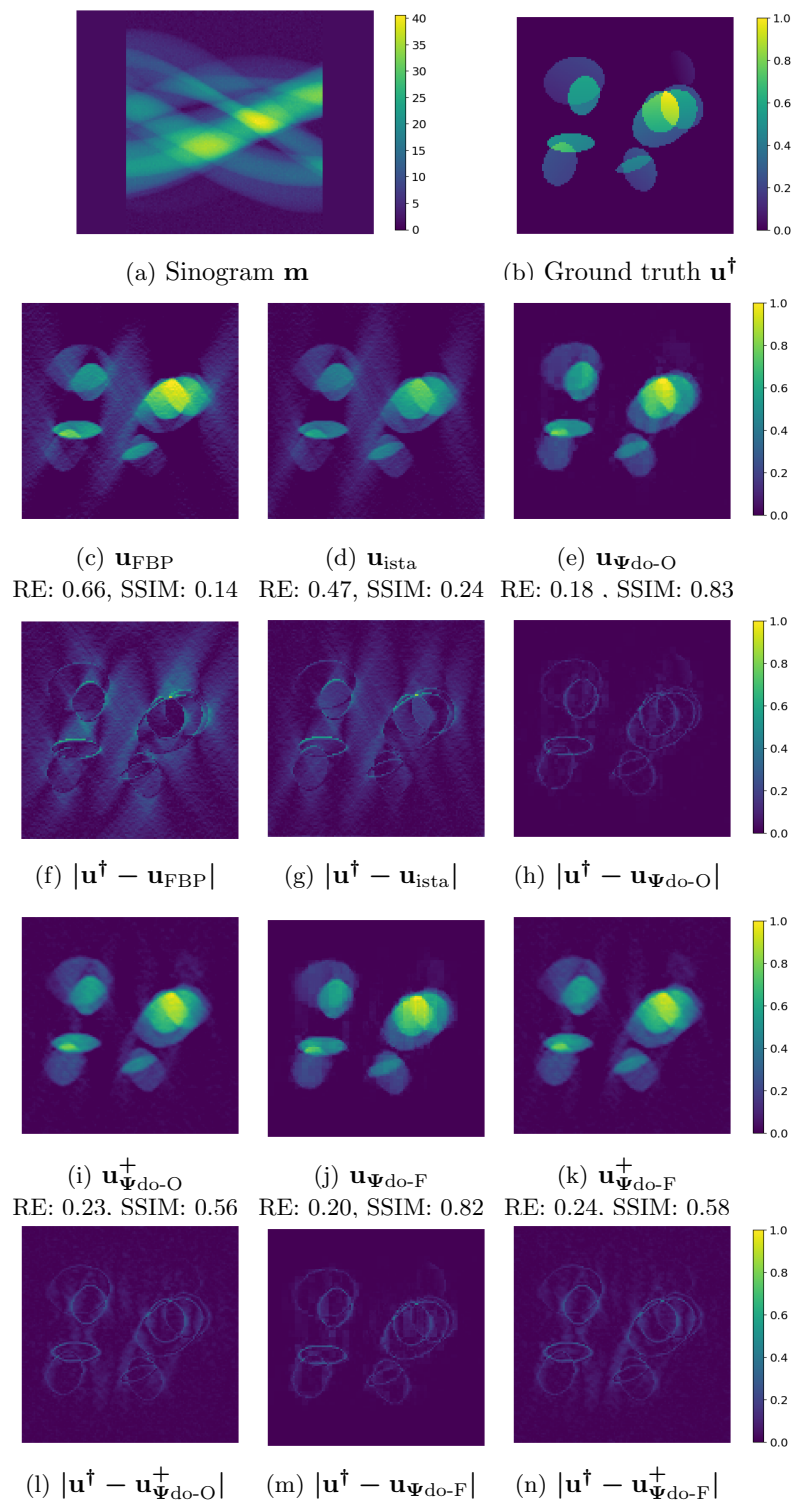
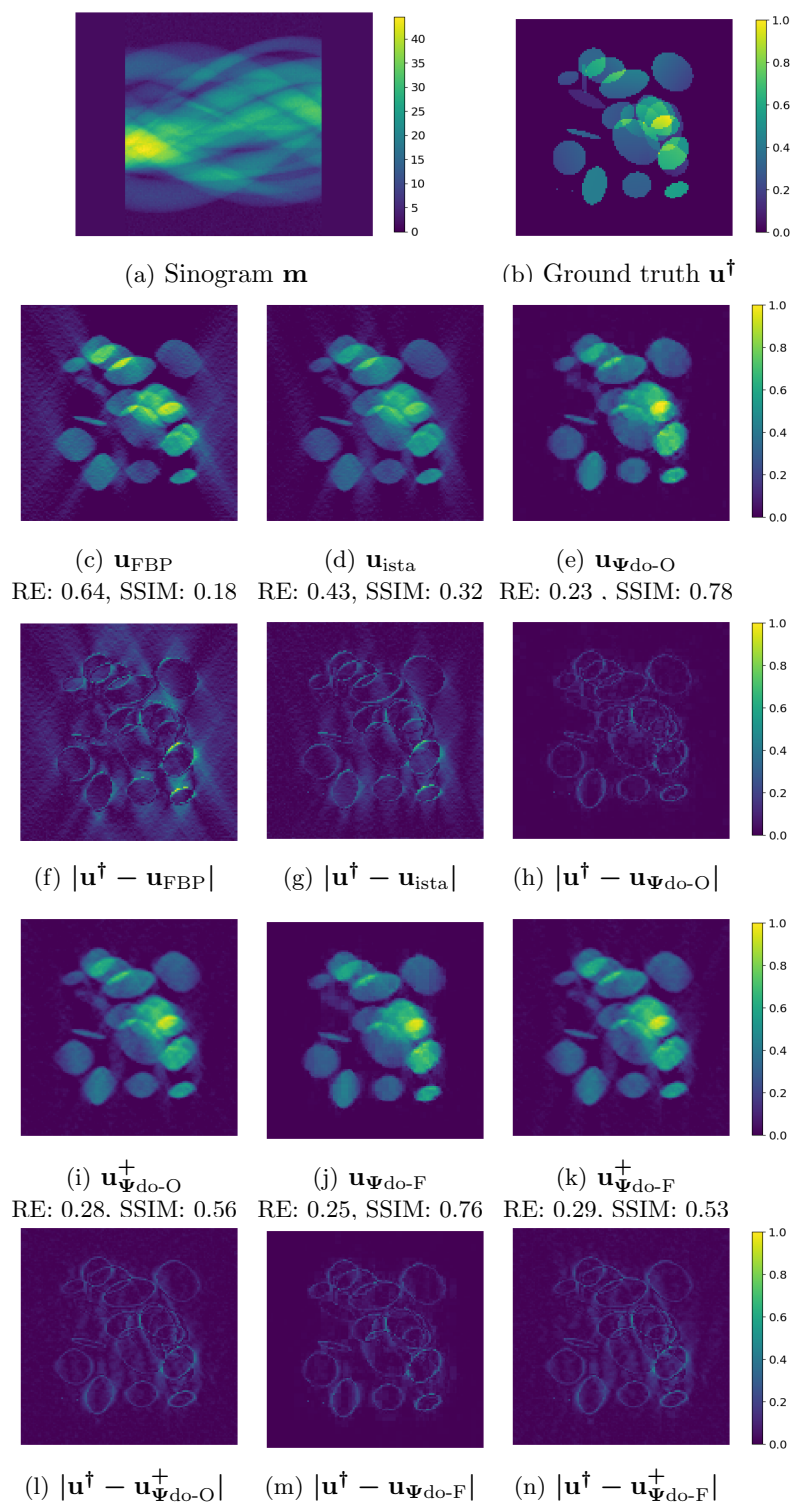


Figure 10. Visualization of the sinogram (or observed image) and corresponding results for one test image.



**Figure 11.** Visualization of the sinogram (or observed image) and corresponding results for one test image.

standard ISTA offers reconstructions of higher quality (cf. Figures 10(d) and 11(d)), however, the streaking artifacts are still noticeable as well as the impurities due to the noise in the measurements. Besides, the ISTA reconstructions are toned down, meaning that for the most part, the intensity of the pixels remain significantly lower than the expected values. With our two models,  $\Psi$ DONet-F and  $\Psi$ DONet-O, whether with positivity constraint on the soft-thresholding parameter or without, it is possible to substantially reduce those artifacts and contrast issues. As can be seen in Figures 10 and 11, our proposed methods lead to undeniably enhanced reconstructions, with a meaningful diminution of the relative error. In particular,  $\Psi$ DONet-O provides slightly better similarity values than  $\Psi$ DONet-F, although both implementations produce comparable results.

In the case where the positivity of the soft-thresholding parameter is enforced, that is for  $\mathbf{u}_{\Psi\text{do-F}}^+$  and  $\mathbf{u}_{\Psi\text{do-O}}^+$ , one can notice that the streaking artifacts, although greatly lessened when compared with the ISTA images, are still present on the reconstructions (cf. Figures 10(i), 10(k), 11(i), and 11(k)). In fact, the SSIM measures are greater than in the ISTA case, but still clearly below the SSIM values obtained with the nonconstrained version of the two models. The latter ( $\mathbf{u}_{\Psi\text{do-F}}$  and  $\mathbf{u}_{\Psi\text{do-O}}$ ) do a noteworthy job in removing the artifacts and sharpening the edges (cf. Figures 10(e), 10(j), 11(e), and 11(j)).

Overall,  $\Psi$ DONet-O without any constraint on the soft-thresholding parameters offers the best results among the compared methods and allows for an optimized implementation of  $\Psi$ DONet.

**7. Conclusions.** In the present paper, we introduced a novel CNN, named  $\Psi$ DONet, inspired by the well-known ISTA and the convolutional nature of certain FIOs and  $\Psi$ DOs, like the limited-angle Radon transform. We proved that the unrolled iterations of ISTA can be interpreted as layers of a CNN, where the downsampling, upsampling, and convolution operations, typically defining a CNN, can be exactly specified by combining the convolutional nature of the limited-angle Radon transform and basic properties defining an orthogonal wavelet system. In addition, we proved that, for a specific choice of the parameters involved,  $\Psi$ DONet recovers standard ISTA or a perturbation of ISTA, up to a bound on the filters coefficients which we estimated in the case of limited-angle Radon transform.

The key feature of the proposed architecture is its potential to learn  $\Psi$ DO-like structures, which makes it suitable to be extended to any convolutional operator which is a FIO or  $\Psi$ DO. Moreover, the analysis carried out on paper allows one to gain understanding and interpretability of the results, which gives insight into a whole class of inverse problems arising from FIO or  $\Psi$ DO and opens it up for fundamental theoretical generalization results.

As a proof of concept, we tested two different implementations of  $\Psi$ DONet on simulated data from limited-angle geometry, generated from the ellipse data set. The improvement, compared to standard ISTA (and classical FBP), is notable and it is promising for further numerical testing which we leave to future work. Additional directions for future numerical testing include larger sizes for images, smaller and sparser visible wedges, additional regularization for the loss function, and reconstructions from real data. Also, it may be beneficial for the reconstruction to introduce more advanced features in the  $\Psi$ DONet architecture, such as skip connections or other residual blocks. Incorporating such elements, while preserving full interpretability of the network, is left for future work.



### Appendix A. Proof of Proposition 2.2.

*Proof.* According to [17, section 3], the following variational source condition is satisfied for every  $w \in \ell^1(\mathbb{N})$ :

$$(A.1) \quad \beta \|w - w^\dagger\|_{\ell^1} \leq \|w\|_{\ell^1} - \|w^\dagger\|_{\ell^1} + C \|AW^*w - AW^*w^\dagger\|_Y.$$

We aim at applying it to  $w = w_{\delta,p,q} \in W_p \subset \ell^1(\mathbb{N})$ . First consider the term  $\|w\|_{\ell^1} - \|w^\dagger\|_{\ell^1}$  in the right-hand side. Since  $w_{\delta,p,q}$  is a solution of (2.5),

$$\begin{aligned} \lambda \|w_{\delta,p,q}\|_{\ell^1} &= (\|A_{p,q}W^*w_{\delta,p,q} - \mathbb{P}_q m\|_Y^2 + \lambda \|w_{\delta,p,q}\|_{\ell^1}) - \|A_{p,q}W^*w_{\delta,p,q} - \mathbb{P}_q m\|_Y^2 \\ &\leq \|A_{p,q}W^*\mathbb{P}_p w^\dagger - \mathbb{P}_q m\|_Y^2 + \lambda \|\mathbb{P}_p w^\dagger\|_{\ell^1} - \|A_{p,q}W^*w_{\delta,p,q} - \mathbb{P}_q m\|_Y^2, \end{aligned}$$

whence

$$\|w_{\delta,p,q}\|_{\ell^1} - \|w^\dagger\|_{\ell^1} \leq \frac{1}{\lambda} \|A_{p,q}W^*\mathbb{P}_p w^\dagger - \mathbb{P}_q m\|_Y^2 - \frac{1}{\lambda} \|A_{p,q}W^*w_{\delta,p,q} - \mathbb{P}_q m\|_Y^2.$$

We can easily check that  $A_{p,q}W^*\mathbb{P}_p = \mathbb{P}_q AW^*\mathbb{P}_p$ ; then, since  $\|\mathbb{P}_q\|_{Y \rightarrow Y} \leq 1$ , denoting by  $Q = \|A_{p,q}W^*w_{\delta,p,q} - \mathbb{P}_q m\|_Y$ , we have

$$\begin{aligned} \|w_{\delta,p,q}\|_{\ell^1} - \|w^\dagger\|_{\ell^1} &\leq \frac{1}{\lambda} \|AW^*\mathbb{P}_p w^\dagger - m\|_Y^2 - \frac{1}{\lambda} Q^2 \\ &\leq \frac{1}{\lambda} \|AW^*(\mathbb{P}_p w^\dagger - w^\dagger)\|_Y^2 + \frac{1}{\lambda} \|AW^*w^\dagger - m\|_Y^2 - \frac{1}{\lambda} Q^2. \end{aligned}$$

In conclusion,

$$(A.2) \quad \|w_{\delta,p,q}\|_{\ell^1} - \|w^\dagger\|_{\ell^1} \leq \frac{1}{\lambda} \|A\|^2 \|w^\dagger - \mathbb{P}_p w^\dagger\|_{\ell^2}^2 + \frac{1}{\lambda} \delta^2 - \frac{1}{\lambda} Q^2.$$

The second term in the right-hand side of (A.1), instead, can be bounded as follows:

$$(A.3) \quad \begin{aligned} \|AW^*(w_{\delta,p,q} - w^\dagger)\|_Y &= \|\mathcal{P}_q AW^*(w_{\delta,p,q} - w^\dagger)\|_Y + \|(I - \mathcal{P}_q)AW^*(w_{\delta,p,q} - w^\dagger)\|_Y \\ &\leq \|A_{p,q}W^*w_{\delta,p,q} - \mathbb{P}_q m\|_Y + \delta + \|(I - \mathbb{P}_q)A\|_{X \rightarrow Y} \|w_{\delta,p,q} - w^\dagger\|_{\ell^1} + \delta \\ &\leq Q + M \|(I - \mathbb{P}_q)A\|_{X \rightarrow Y} + \delta, \end{aligned}$$

where the positive constant  $M$  depends on  $\|w^\dagger\|_{\ell^2}$ . In order to get an estimate for  $Q = \|A_{p,q}W^*w_{\delta,p,q} - \mathbb{P}_q m\|_Y$ , we use (A.1): since  $0 \leq \beta \|w_{\delta,p,q} - w^\dagger\|_{\ell^1}$ , using (A.2) and (A.3) we have

$$0 \leq \frac{1}{\lambda} \|A\| \|w^\dagger - \mathbb{P}_p w^\dagger\|_{\ell^2}^2 + \frac{1}{\lambda} \delta^2 - \frac{1}{\lambda} Q^2 + Q + M \|(I - \mathbb{P}_q)A\|_{X \rightarrow Y} + \delta.$$

By solving this second-order inequality we get

$$(A.4) \quad \begin{aligned} Q &\leq \frac{\lambda}{2} + \frac{\lambda}{2} \left( 1 + \frac{4}{\lambda^2} \|A\|^2 \|w^\dagger - \mathbb{P}_p w^\dagger\|_{\ell^2}^2 \frac{4}{\lambda} \delta^2 + \frac{4M}{\lambda} \|(I - \mathbb{P}_q)A\|_{X \rightarrow Y} \frac{4}{\lambda} \delta \right)^{\frac{1}{2}} \\ &\leq \lambda + \delta + \|A\| \|w^\dagger - \mathbb{P}_p w^\dagger\|_{\ell^2} + M \|(I - \mathbb{P}_q)A\|_{X \rightarrow Y}. \end{aligned}$$

Combining (A.1), (A.2), (A.3), and (A.4) we easily conclude the proof. ■

**Appendix B. Proof of Proposition 2.5.**

*Proof.* Consider the sequence  $e_n = \|\mathbf{w}_\rho^{(n+1)} - \mathbf{w}^{(n+1)}\|_{\ell^2}$ . Thanks to the nonexpansivity of the operator  $\mathcal{S}_{\lambda/L}$ , it holds that

$$(B.1) \quad e_0 = \|\mathcal{T}_Z(\mathbf{w}^{(0)}) - \mathcal{T}(\mathbf{w}^{(0)})\|_{\ell^2} \leq \frac{1}{L} \|\mathbf{W} \mathbf{A}_{p,q}^* \mathbf{A}_{p,q} \mathbf{W}^* - \mathbf{Z}\| \|\mathbf{w}^{(0)}\|_{\ell^2} \leq \frac{1}{L} \rho \|\mathbf{w}^{(0)}\|_{\ell^2}.$$

Analogously, for  $n \geq 1$ ,

$$(B.2) \quad \begin{aligned} e_n &\leq \left\| \mathbf{I} - \frac{1}{L} \mathbf{Z} \right\| \|\mathbf{w}^{(n)} - \mathbf{w}_\rho^{(n)}\|_{\ell^2} + \frac{1}{L} \|\mathbf{W} \mathbf{A}_{p,q}^* \mathbf{A}_{p,q} \mathbf{W}^* - \mathbf{Z}\| \|\mathbf{w}^{(n)}\|_{\ell^2} \\ &\leq \left\| \mathbf{I} - \frac{1}{L} \mathbf{Z} \right\| e_{n-1} + \frac{1}{L} \rho \|\mathbf{w}^{(n)}\|_{\ell^2}. \end{aligned}$$

Since  $L \geq \|\mathbf{W} \mathbf{A}_{p,q}^* \mathbf{A}_{p,q} \mathbf{W}^*\|$ , then

$$\left\| \mathbf{I} - \frac{1}{L} \mathbf{Z} \right\| \leq \left\| \mathbf{I} - \frac{1}{L} \mathbf{W} \mathbf{A}_{p,q}^* \mathbf{A}_{p,q} \mathbf{W}^* \right\| + \frac{1}{L} \|\mathbf{W} \mathbf{A}_{p,q}^* \mathbf{A}_{p,q} \mathbf{W}^* - \mathbf{Z}\| \leq 1 + \frac{1}{L} \rho.$$

Moreover, since the sequence  $\{\mathbf{w}^{(n)}\}$  is convergent, then it is also bounded: let, e.g.,  $\|\mathbf{w}^{(n)}\|_{\ell^2} \leq M$ . As a consequence of (B.1), (B.2),

$$e_N \leq \sum_{n=0}^N \left(1 + \frac{1}{L} \rho\right)^{N-n} \frac{1}{L} \rho \|\mathbf{w}^{(n)}\|_{\ell^2} \leq M \left( \left(1 + \frac{1}{L} \rho\right)^{N+1} - 1 \right).$$

Let now  $N \geq N_0$  and  $\rho N \leq \eta_0$ : then, with a constant  $c = c(N_0, \eta_0)$ , it holds that

$$\|\mathbf{w}_\rho^{(N)} - \mathbf{w}^{(N)}\|_{\ell^2} = e_N \leq M(e^{\frac{1}{L} \rho N} - 1) \leq c \frac{M}{L} \rho N.$$

Combining this result with (2.8), we can guarantee that

$$\|\mathbf{w}_\rho^{(N)} - \mathbf{w}_{\delta,p,q}\|_{\ell^2} \leq \|\mathbf{w}_\rho^{(N)} - \mathbf{w}^{(N)}\|_{\ell^2} + \|\mathbf{w}^{(N)} - \mathbf{w}_{\delta,p,q}\|_{\ell^2} \leq c_3 a^N + \tilde{c}_4 \rho N,$$

which proves (2.11). To obtain (2.12), simply substitute  $N = \log_a \delta$  and  $\rho = \frac{\delta}{N}$  and consider  $c_4 = c_3 + \tilde{c}_4$ . ■

**REFERENCES**

- [1] J. ADLER, H. KOHR, AND O. ÖKTEM, *Operator discretization library (ODL)*, 2017, [https://zenodo.org/record/556409#.YHw3rz\\_ONPY](https://zenodo.org/record/556409#.YHw3rz_ONPY).
- [2] J. ADLER AND O. ÖKTEM, *Solving ill-posed inverse problems using iterative deep neural networks*, *Inverse Problems*, 33 (2017), 124007.
- [3] J. ADLER AND O. ÖKTEM, *Learned primal-dual reconstruction*, *IEEE Trans. Med. Imaging*, 37 (2018), pp. 1322–1332.
- [4] S. ARRIDGE, P. MAASS, O. ÖKTEM, AND C. SCHÖNLIEB, *Solving inverse problems using data-driven models*, *Acta Numer.*, 28 (2019), pp. 1–174.

- [5] G. BEYLKIN, R. COIFMAN, AND V. ROKHLIN, *Fast wavelet transforms and numerical algorithms I*, Comm. Pure Appl. Math., 44 (1991), pp. 141–183.
- [6] J. BOLTE, T. NGUYEN, J. PEYPOUQUET, AND B. SUTER, *From error bounds to the complexity of first-order descent methods for convex functions*, Math. Program., 165 (2017), pp. 471–507.
- [7] L. BORG, J. FRIKEL, J. S. JORGENSEN, AND E. T. QUINTO, *Analyzing reconstruction artifacts from arbitrary incomplete X-ray CT data*, SIAM J. Imaging Sci., 11 (2018), pp. 2786–2814.
- [8] K. BREDIES AND D. LORENZ, *Linear convergence of iterative soft-thresholding*, J. Fourier Anal. Appl., 14 (2008), pp. 813–837.
- [9] T. A. BUBBA, G. KUTYNIOK, M. LASSAS, M. MÄRZ, W. SAMEK, S. SILTANEN, AND V. SRINIVASAN, *Learning the invisible: A hybrid deep learning-shearlet framework for limited angle computed tomography*, Inverse Problems, 35 (2019), 064002.
- [10] W. DAHMEN, *Wavelet and multiscale methods for operator equations*, Acta Numer., 6 (1997), pp. 55–228.
- [11] W. DAHMEN, S. PRÖSSDORF, AND R. SCHNEIDER, *Wavelet approximation methods for pseudodifferential equations: I Stability and convergence*, Math. Z., 215 (1994), pp. 583–620.
- [12] I. DAUBECHIES, M. DEFRISE, AND C. DE MOL, *An iterative thresholding algorithm for linear inverse problems with a sparsity constraint*, Comm. Pure Appl. Math., 57 (2004), pp. 1413–1457.
- [13] M. E. DAVISON, *The ill-conditioned nature of the limited angle tomography problem*, SIAM J. Appl. Math., 43 (1983), pp. 428–448.
- [14] M. DE HOOP, M. LASSAS, AND C. WONG, *Deep Learning Architectures for Nonlinear Operator Functions and Nonlinear Inverse Problems*, preprint, <https://arxiv.org/abs/1912.11090> (2019).
- [15] D. FANELLI AND O. ÖKTEM, *Electron tomography: A short overview with an emphasis on the absorption potential model for the forward problem*, Inverse Problems, 24 (2008), 013001.
- [16] D. FINCH, I.-R. LAN, AND G. UHLMANN, *Microlocal analysis of the x-ray transform with sources on a curve*, in Inside Out: Inverse Problems and Applications, G. Uhlmann, ed., Cambridge University Press, Cambridge, 2003, pp. 193–218.
- [17] J. FLEMMING AND D. GERTH, *Injectivity and weak\*-to-weak continuity suffice for convergence rates in  $\ell_1$ -regularization*, J. Inverse Ill-Posed Probl., 26 (2018), pp. 85–94.
- [18] J. FRIKEL AND E. QUINTO, *Characterization and reduction of artifacts in limited angle tomography*, Inverse Problems, 29 (2013), 125007.
- [19] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep Learning*, MIT Press, Cambridge, MA, 2016.
- [20] A. GREENLEAF AND G. UHLMANN, *Estimates for singular radon transforms and pseudodifferential operators with singular symbols*, J. Funct. Anal., 89 (1990), pp. 202–232.
- [21] A. GREENLEAF AND G. UHLMANN, *Nonlocal inversion formulas for the x-ray transform*, Duke Math. J., 58 (1989), pp. 205–240.
- [22] K. GREGOR AND Y. LECUN, *Learning fast approximations of sparse coding*, in Proc. 27th ICML, International Machine Learning Society, Madison, WI, 2010, pp. 399–406.
- [23] K. HAUG AND M. LOHNE, *TF-Wavelets*, <https://github.com/UiO-CS/tf-wavelets> (2019).
- [24] K. HEISKANEN, H. RHIM, AND P. MONTEIRO, *Computer simulations of limited angle tomography of reinforced concrete*, Cement Concrete Res., 21 (1991), pp. 625–634.
- [25] J. R. HERSHEY, J. L. ROUX, AND F. WENINGER, *Deep Unfolding: Model-Based Inspiration of Novel Deep Architectures*, preprint <https://arxiv.org/abs/1409.2574> (2014).
- [26] L. HÖRMANDER, *The Analysis of Linear Partial Differential Operators*, Vol. III, Springer, Berlin, 1985.
- [27] L. HÖRMANDER, *The Analysis of Linear Partial Differential Operators*, Vol. IV, Springer, Berlin, 1985.
- [28] K. JIN, M. MCCANN, E. FROUSTEY, AND M. UNSER, *Deep convolutional neural network for inverse problems in imaging*, IEEE Trans. Image Process., 26 (2017), pp. 4509–4522.
- [29] A. I. KATSEVICH, *Local tomography for the generalized Radon transform*, SIAM J. Appl. Math., 57 (1997), pp. 1128–1162.
- [30] A. KATSEVICH, *Local tomography for the limited-angle problem*, J. Math. Anal. Appl., 213 (1997), pp. 160–182.
- [31] D. KINGMA AND J. BA, *Adam: A Method for Stochastic Optimization*, preprint, arXiv:1412.6980 (2015).
- [32] V. KOLEHMAINEN, S. SILTANEN, S. JÄRVENPÄÄ, J. P. KAIPIO, P. KOISTINEN, M. LASSAS, J. PIRTTILÄ, AND E. SOMERSALO, *Statistical inversion for medical x-ray tomography with few radiographs: II. Application to dental radiology*, Phys. Med. Biol., 48 (2003), 1465.

- [33] V. KRISHNAN AND E. QUINTO, *Microlocal analysis in tomography*, in Handbook of Mathematical Methods in Imaging, Springer, New York, 2015, pp. 847–902.
- [34] G. LEE, R. GOMMERS, F. WASELEWSKI, K. WOHLFAHRT, AND A. O’LEARY, *PyWavelets: A Python package for wavelet analysis*, J. Open Source Softw., 4 (2019), 1237.
- [35] S. LUNZ, O. ÖKTEM, AND C.-B. SCHÖNLIEB, *Adversarial regularizers in inverse problems*, in Proc. 32nd NIPS, Curran Associates, Red Hook, NY, 2018, pp. 8507–8516.
- [36] S. MALLAT, *A wavelet tour of signal processing*, Elsevier, San Diego, CA, 1999.
- [37] J. L. MUELLER AND S. SILTANEN, *Linear and Nonlinear Inverse Problems with Practical Applications*, Comput. Sci. Eng. 10, SIAM, Philadelphia, 2012.
- [38] F. NATTERER, *The mathematics of computerized tomography*, Classics Appl. Math. 32, SIAM, Philadelphia, 2001.
- [39] F. NATTERER AND F. WÜBBELING, *Mathematical methods in image reconstruction*, Math. Model. Comput., SIAM, Philadelphia, 2001.
- [40] L. V. NGUYEN, *How strong are streak artifacts in limited angle computed tomography?*, Inverse Problems, 31 (2015), 055003.
- [41] L. V. NGUYEN, *On the strength of streak artifacts in filtered back-projection reconstructions for limited angle weighted X-ray transform*, J. Fourier Anal. Appl., 23 (2017), pp. 712–728.
- [42] C. OLIVER, *Synthetic-aperture radar imaging*, J. Phys. D, 22 (1989), pp. 871–890.
- [43] E. T. QUINTO, *Singularities of the X-Ray transform and limited data tomography in  $\mathbb{R}^2$  and  $\mathbb{R}^3$* , SIAM J. Math. Anal., 24 (1993), pp. 1215–1225.
- [44] E. QUINTO, *An introduction to X-ray tomography and Radon transforms*, in The Radon Transform, Inverse Problems, and Tomography, Proc. Sympos. Appl. Math. 63, American Mathematical Society, Providence, RI, 2006, pp. 1–23.
- [45] A. RAMM AND A. KATSEVICH, *Inversion of incomplete Radon transform*, Appl. Math. Lett., 5 (1992), pp. 41–45.
- [46] A. G. RAMM AND A. I. KATSEVICH, *The Radon Transform and Local Tomography*, CRC Press, Boca Raton, FL, 1996.
- [47] R. REISENHOFER, S. BOSSE, G. KUTYNIOK, AND T. WIEGAND, *A Haar wavelet-based perceptual similarity index for image quality assessment*, Signal Process. Image Comm., 61 (2018), pp. 33–43.
- [48] D. RUMELHART, G. HINTON, AND R. WILLIAMS, *Learning representations by back-propagating errors*, Nature, 323 (1986), pp. 533–536.
- [49] R. TOVEY, M. BENNING, C. BRUNE, M. J. LAGERWERF, S. M. COLLINS, R. K. LEARY, P. A. MIDGLEY, AND C.-B. SCHÖNLIEB, *Directional sinogram inpainting for limited angle tomography*, Inverse Problems, 35 (2019), 024004.
- [50] G. UHLMANN AND A. VASY, *The inverse problem for the local geodesic ray transform*, Invent. Math., 205 (2016), pp. 83–120.
- [51] W. VAN AARLE, W. J. PALENSTIJN, J. CANT, E. JANSSENS, F. BLEICHRODT, A. DABRAVOLSKI, J. D. BEENHOUWER, J. BATENBURG, AND J. SIJBERS, *Fast and flexible X-ray tomography using the ASTRA toolbox*, Opt. Express, 24 (2016), pp. 25129–25147.
- [52] S. VAN DER WALT, J. L. SCHÖNBERGER, J. NUNEZ-IGLESIAS, F. BOULOGNE, J. D. WARNER, N. YAGER, E. GOUILLART, AND T. YU, *scikit-image: Image processing in Python*, Peer J. 2004 (2014) e453.
- [53] Z. WANG, A. BOVIK, H. SHEIKH, AND E. P. SIMONCELLI, *Image quality assessment: From error visibility to structural similarity*, IEEE Trans. Image Process., 13 (2004), pp. 600–612.
- [54] J. ZHANG AND B. GHANEM, *ISTA-Net: Interpretable optimization-inspired deep network for image compressive sensing*, in Proc. IEEE Conf. Computer Vision and Pattern Recognition, IEEE Computer Society, Los Alamitos, CA, 2018, pp. 1828–1837.
- [55] K. ZHANG, L. V. GOOL, AND R. TIMOFTE, *Deep unfolding network for image super-resolution*, in Proc. IEEE Computer Society Conf. Computer Vision Pattern Recognition, IEEE Computer Society, Los Alamitos, CA, 2020, pp. 3217–3226.
- [56] Y. ZHANG, H. P. CHAN, B. SAHNER, J. WEI, M. GOODSITT, L. M. HADJIISKI, J. GE, AND C. ZHOU, *A comparative study of limited-angle cone-beam reconstruction methods for breast tomosynthesis*, Med. Phys., 33 (2006), pp. 3781–3795.